

# MEAV: Model Editing with Alignment Vectors for inference time LLM alignment in single and multidomain preference spectrum

Anonymous ACL submission

## Abstract

Aligning LLMs to nuanced preference levels requires adequate flexibility and control, which can be a resource-intensive and time-consuming procedure. Existing training-time alignment methods require full re-training when a change is needed and inference-time ones typically require access to the reward model at each inference step. We introduce **MEAV**, an inference-time model-editing-based LLM alignment method that learns encoded representations of preference dimensions, called *Alignment Vectors* (AV). These representations enable dynamic adjusting of the model behavior during inference through simple linear operations. Here, we focus on three gradual response levels across three specialized domains: medical, legal, and financial, exemplifying its practical potential. We introduce adjustable preference knobs during inference, allowing users to tailor their LLM outputs while reducing the inference cost by half compared to the prompt engineering approach. Additionally, AVs are transferable across different fine-tuning stages of the same model, demonstrating flexibility. AVs also facilitate multidomain, diverse preference alignment, making the process 12x faster than the retraining approach.

## 1 Introduction

Aligning LLMs is crucial for adapting them to meet human preferences. Standard training-time alignment methods, such as RLHF (Ouyang et al., 2022) and DPO (Rafailov et al., 2024), are conducted during model training. However, making nuanced preference adjustments during inference with these approaches necessitates retraining, which requires substantial amounts of time, preference data and computational resources. Inference-time LLM alignment, by contrast, delays the alignment process until inference (Wang et al., 2024). While preference alignment can be achieved through training-time methods or targeted prompting, fine-

grained control over preferences at inference remains largely unexplored in current State-of-the-Art (SOTA) works (Sahoo et al., 2024; Guo et al., 2024). This research introduces an inference-time model editing technique via *Alignment Vectors* (AV), offering users dynamic preference adjustments without additional computational overhead.

Due to their extensive capabilities, LLMs are now employed in different fields, but the diverse needs of a broad customer base require that LLM outputs be carefully refined. For instance, while a healthcare provider might need detailed medical responses for professional use, a public health forum may prefer more generalized information to avoid misinterpretation. Although prompt engineering can temporarily address these needs, it becomes costly when scaled (Li et al., 2023).

Furthermore, managing multiple alignment objectives can be complex. Consider an insurance company that needs expert legal responses, generic financial answers, and to avoid medical responses; balancing these demands poses a significant challenge. A joint training with targeted preference levels can resolve the problem, however, it lacks flexibility, and training can be resource intensive. Hence, at present, there is no work that addresses such preference flexibility in the inference time. Thus, developing flexible, inference-time adjustable model alignment to manage costs and maintain efficiency in the long term remains a major research gap.

Preference dimensions like helpfulness, harmlessness, and honesty are well-studied, with some work exploring their controllability via numerical levels (Bai et al., 2022; Ji et al., 2024; Guo et al., 2024). However, specialized dimensions offer finer granularity, enabling better control during inference. To enhance preference tunability, we focus on proficiency levels in specialized domains while also demonstrating tunability in a general domain, such as safety. Since existing literature lacks

084 domain-specific preference alignment datasets, we  
085 generate synthetic Query-Response pairs by deriv-  
086 ing queries from the PersonaHub dataset (Chan  
087 et al., 2024) and augmenting them with novel per-  
088 sonas created via LLM-generated prompts.

089 In addition, to achieve inference time preference  
090 tunability, we propose a simple technique called  
091 **Model Editing via Alignment Vector (MEAV)**,  
092 which is based on the concept of *Task Arithmetic* (Il-  
093 harco et al., 2023). AVs can be obtained by directly  
094 subtracting the base model parameters from the  
095 aligned model, and can be added in the inference  
096 time. Hence, our first research question (**RQ1**) Are  
097 alignment vectors valid representation of the pref-  
098 erence dimensions? To address this question, we  
099 systematically integrate the alignment vector into  
100 the base model with varying weights, both positive  
101 and negative, and analyze the resulting changes in  
102 model behavior. Our second research question is  
103 posed as (**RQ2**) Can we calibrate different align-  
104 ment vectors to achieve diverse multi-domain pref-  
105 erence? We address RQ2 through different domain-  
106 specific AV-integration strategy.

107 The key contribution of this work are:

- 108 • We frame LLM alignment in single and multi-  
109 ple domains as a model editing problem and  
110 introduce an inference-time tunable mecha-  
111 nism, which allows flexible adjustment of  
112 generation output along the preference dimen-  
113 sion.
- 114 • We generate a synthetic dataset with a total of  
115 38k queries, each paired with responses cate-  
116 gorized into three levels of specialized subject  
117 matter proficiency across three specialized do-  
118 mains: Medical, Financial, and Legal. The  
119 dataset will be available through this link.
- 120 • By adjusting the merging coefficients, we  
121 achieve diverse, multidomain behaviors effi-  
122 ciently, saving time and resources. Unlike  
123 joint training, which requires  $p^D$  adjustments  
124 for  $D$  domains and  $p$  preference levels, our  
125 method only requires  $D$  training runs, reduc-  
126 ing resource usage by a factor of  $p^D/D$ .

## 127 2 Related Works

128 While prompt engineering techniques are effective  
129 in aligning LLM responses to user queries during  
130 inference time, they incur high inference costs and  
131 rely heavily on user expertise on prompting (Rad-

ford et al., 2019; Meskó, 2023; Oppenlaender et al.,  
2023).

132 Li et al. introduced an inference-time techni-  
133 que that identifies a sparse set of attention heads  
134 for a target task and shifts their activation along  
135 task-correlated directions during inference time (Li  
136 et al., 2024). A similar approach was explored to  
137 learning Safety Related Vectors, to steer harmful  
138 model outputs towards safer alternatives (Wang  
139 et al., 2024). However, these methods were tar-  
140 get domain-specific and not controllable. Huang  
141 et al. introduced *DeAl*, an alignment method that  
142 treats alignment as a heuristic-guided search pro-  
143 cess (Huang et al., 2024). Liu et al. studied regu-  
144 larization strength between aligned and unaligned  
145 models to have control over generation (Liu et al.,  
146 2024). Although closely related to our work, their  
147 method lacks clarity on whether fine-grained pref-  
148 erence levels can be achieved. Researchers controlled  
149 attributes of generated contents by adding  
150 control token in the prompt (Guo et al., 2024; Dong  
151 et al., 2023). Despite its effectiveness, this method  
152 requires training LLMs with a particular data for-  
153 mat, which restricts the flexibility of control during  
154 inference.

155 Rame et al.’s work is closely related to our multi-  
156 domain preference alignment where they merge  
157 multiple fully fine-tuned models (each trained on  
158 a different reward) into a single “soup,” typically  
159 yielding a fixed blend (Rame et al., 2023). In con-  
160 trast, MEAV obtains a single difference vector per  
161 domain or preference dimension and adds it to the  
162 base model at inference time, allowing a smooth,  
163 continuous controllability for alignment. As a re-  
164 sult, Rewarded Soups requires training multiple  
165 models upfront, whereas MEAV simply reuses one  
166 base model with tunable additive vectors. As for  
167 the comparison, it is not clear how the controlla-  
168 bility objective (which is our main focus) can be  
169 achieved in the rewarded soup model. Similarly,  
170 while Jang et al. address personalized preference  
171 alignment and post-hoc merging, our approach pro-  
172 vides a unique capability: preference level adjust-  
173 ment (Jang et al., 2023). In summary, while pa-  
174 rameter merging techniques relate to our goal of  
175 multi-domain alignment, MEAV is distinct in en-  
176 abling real-time, granular control over preference  
177 strength, an ability not offered by existing merg-  
178 ing-based methods.

179 Yang et al. used a multi-objective training ap-  
180 proach that encodes multiple reward signals within  
181 a single model, adjusting preferences via special-  
182  
183

ized prompts or latent context (Yang et al., 2024). In contrast, MEAV learns a distinct additive vector per domain or preference dimension, which can be applied to a base model at inference time for continuous, fine-grained control, without re-training. Similarly, Yu et al. absorbed additional capabilities from models with the same architecture by merging their parameters, effectively accruing new skills as a “free lunch” (Yu et al., 2024). While MEAV also uses a form of parameter combination (subtracting and adding AVs), MEAV targets preference alignment at inference time, not transferring new tasks or knowledge.

### 3 Methodology

#### 3.1 Obtaining Alignment Vector

To obtain the AVs, we first perform alignment through the DPO algorithm, using an ‘ipo’ loss function to create a domain-specific aligned model (Rafailov et al., 2024; Azar et al., 2024). We get AVs by subtracting the weights of an unaligned model from the weights of the same model after alignment on a task. If  $\theta_{aligned}$  denotes the model parameter after aligning on a preference dimension, then the AV can be obtained by the following:

$$\theta_{AV} = \theta_{aligned} - \theta_{unaligned} \quad (1)$$

#### 3.2 Single Domain Alignment

To enable preference tunability across different domains, we perform a weighted integration of the AVs into the base (or unaligned) model, where the weights can be both positive and negative. We hypothesize that this gradual integration will result in a corresponding gradual increase or decrease in the model’s proficiency. This process is governed by the following equation.

$$\theta_{aligned} = \theta_{unaligned} + \lambda * \theta_{AV} \quad (2)$$

By adjusting the value of  $\lambda$ , we aim to control the proficiency of the model’s generated responses. Assuming when  $\lambda = 0$ , the model remains unaltered and functions as the base, unaligned model. If the  $\theta_{AV}$  encodes the expert behavior in a certain domain, as  $\lambda$  increases towards 1, the model becomes increasingly aligned, achieving full proficiency at  $\lambda = 1$ .

We further hypothesize that when  $\lambda$  takes on negative values, the model’s behavior tends to reverse the preference ranking. For instance, if the base model typically generates generic responses

and the aligned model is designed for expert-level responses, moving  $\lambda$  in the negative direction will shift the model towards avoidance behavior. Therefore, to control the proficiency of the responses, adjusting  $\lambda$  is sufficient, eliminating the need to train the model with a new preference configuration.

#### 3.3 Multidomain Alignment

When dealing with multiple domains simultaneously, the interaction between these domains can present a significant challenge. While individual preference vector encodes domain-specific attributes, they also embed proficiency levels which can easily generalize and negatively affect multidomain diverse behavior. This complexity can make it difficult to integrate multiple domains effectively.

Our goal is to achieve a diverse multidomain preference, which we approach by using the following equation:

$$\theta_{multidom\_aligned} = \theta_{unaligned} + \alpha\theta_{AV\_dom1} + \beta\theta_{AV\_dom2} + \gamma\theta_{AV\_dom3} \quad (3)$$

In this equation,  $\alpha$ ,  $\beta$  and  $\gamma$  represent the integration coefficients for the domains in question, respectively. By identifying different sets of these coefficients, we aim to achieve varying levels of preference across the three domains.

### 4 Synthesizing Specialized Preference Data

To gather data for preference tuning on response proficiency levels, we employ two methods to collect queries: “PersonaHub” (Chan et al., 2024) and “CreatePersona.” Figure 1 provides a detailed overview of the process. Notably, all generated persona, queries, responses, and the prompts used are in English.

#### 4.1 Query Generation

We initiate the generation with a hierarchical process called “CreatePersona.” We begin by randomly generating a few persona-query pairs by prompting Claude-3-Sonnet (Anthropic, 2024). To preserve diversity, we limit the initial set to five pairs, as we found generating too many at the outset reduces variation. From each initial persona, we recursively generate additional persona-query pairs that are relevant to the root persona. We randomize this process three times.

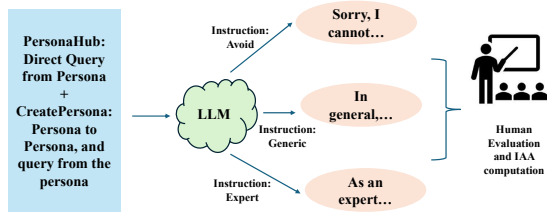


Figure 1: The process of data collection. Personas are sourced from both the PersonaHub dataset and the CreatePersona method. These personas are then fed to an LLM to generate queries. The LLM is prompted with specific instructions to produce responses across three proficiency levels. Following this, human evaluation is conducted to ensure the accuracy and quality of the generated response levels.

To further diversify the dataset, we supplement our generated personas by randomly sampling an equal number from the PersonaHub dataset (Chan et al., 2024), licensed as cc-by-nc-sa-4.0. Using these selected personas, we prompt Claude-3-Sonnet (Anthropic, 2024) to generate specialized domain queries.

We chose Claude-3-Sonnet over GPT-4 for two main reasons: First, Claude-3-Sonnet has consistently demonstrated performance on par with GPT-4, often ranking among the best foundational models. Second, we opted to use GPT-4 as an independent evaluator and sought to mitigate the known bias where evaluators tend to favor their own outputs over those generated by other models (Zheng et al., 2024; Anthropic, 2024).

After a thorough clean-up, involving truncation, and reformatting, we obtained 13,000 personas for the medical domain, 12,374 personas for the financial domain, and 12,867 personas for the legal domain. Each persona is accompanied by queries pertinent to their respective specialized domains.

## 4.2 Response Generation

We generate the response from the queries into three distinct levels: avoidance of response (Avd), generic response (Gen), and expert response (Exp). Detailed instructions are provided to the LLM to facilitate the generation of these responses (see Appendix C). Furthermore, we observe a progressive increase in response length from the avoidance level to the expert level. To mitigate potential bias associated with response length, we instructed the LLM to produce responses of random lengths.

Across medical, financial, and legal domains, we evaluated factual correctness using a 50-sample

experiment per domain, where GPT-4o (OpenAI, 2025a) and GPT-o3-mini (OpenAI, 2025b) independently scored expert responses on a 10-point correctness scale. Both models consistently assigned high marks: GPT-4o with averages of 9.10 (medical), 8.62 (financial), and 8.62 (legal), and GPT-o3-mini with even stronger scores of 9.74, 9.42, and 9.66 in the same domains, demonstrating robust factual alignment of the generated expert responses without requiring domain-specific adjustments.

## 4.3 Human Evaluation of multi-level response generation

To evaluate the quality of the generated responses, we conduct a small experiment involving three annotators, and compute the Inter-Annotator Agreement (IAA). Each annotator is asked to categorize a set of LLM-generated responses into one of three categories: Avd, Gen, and Exp. We provide the annotators with clear definitions of these categories. Each annotator reviews 30 queries along with their three-level responses, with at least 15 examples shared between every pair of annotators. This allows us to compute the average Cohen’s kappa score, which is found to be 0.84 (Cohen, 1960), indicating substantial agreement among the annotators.

We also calculate the average annotation agreement for each annotator with the LLM generation. Responses generated with the Avoidance instruction have the fewest disagreements or misclassifications. However, some Gen and Exp responses are occasionally misclassified from one another. We observe that certain responses, although aligned with the expert spectrum, are misidentified as generic due to their tone, and vice versa. Additionally, a few avoidance responses provide basic information, leading to their misclassification as Gen responses. These findings suggest that the levels may represent a continuous spectrum rather than distinct categories, highlighting the need for further research to more precisely define these proficiency levels.

## 5 Experiments

### 5.1 Evaluation Metric

To assess the performance after alignment, we use a metric called *preference accuracy* (pref. acc). This metric reports the accuracy at each alignment level. To calculate it, we first compute the token-level

mean log-probability ( $MLP$ ) for each of the three response levels across all queries for the aligned model. Then, for each sample in the validation set, we determine which alignment level has the highest log-probability. For example, in proficiency level alignment, it can be among Exp, Gen, and Avd. Finally, we report the percentage of samples where each alignment level had the highest log-probability in the validation set. A higher preference accuracy in an alignment spectrum indicate the dominant behavior of that level.

To illustrate, for a query  $q \in Q$ , the mean log-probability for response  $r \in R$ , where  $R$  can be different alignment levels, is computed for model  $M_\lambda$  as:

$$MLP(r, q, M_\lambda) = \frac{1}{T_r(q)} \sum_{i=1}^{T_r(q)} \log P(t_i | \text{ctx}, M_\lambda) \quad (4)$$

where  $T_r(q)$  is the response length,  $t_i$  is the  $i^{\text{th}}$  token and  $\text{ctx}$  is the previously processed context. The preferred alignment level is:

$$r^*(q) = \arg \max_{r \in R} MLP(r, q, M_\lambda).$$

The preference accuracy for level  $r$  is:

$$Pref. Acc(r) = \frac{1}{|Q|} \sum_{q \in Q} \mathbf{1}[r^*(q) = r],$$

where  $\mathbf{1}[r^*(q) = r]$  is the indicator function. Higher  $Pref. Acc(r)$  indicates the dominant behavior of the preference alignment level  $r$ . A similar approach was also used in pairwise preference accuracy computation in (Stiennon et al., 2020).

Additionally, we use an auxiliary metric as ‘‘GPT-4 judged generation accuracy’’, where we generate the responses from queries in a sample, and ask GPT-4 to annotate it as one of the three levels (Zheng et al., 2024). After that, we simply report the percentage of each annotated alignment level.

## 5.2 Baseline Approaches

Since existing model-editing methods lack inference-time controlled alignment, we use ‘prompting’ as a baseline, instructing the LLM to generate responses at predefined proficiency levels. Unlike model editing, this enables discrete levels rather than a spectrum. Our second baseline, ‘Joint Training,’ combines multidomain data to align responses across proficiency levels, offering

insights despite being a training-time method. We also report the model’s ‘default’ performance, where queries are prompted without additional instructions or edits.

## 5.3 Model and Training Configuration

We define three main preference levels: ‘‘expert,’’ ‘‘generic,’’ and ‘‘avoidance’’ for specialized domain proficiency and use DPO training with a fixed beta of 0.1, where ‘‘expert’’ is preferred over ‘‘generic,’’ and ‘‘generic’’ over ‘‘avoidance.’’ To demonstrate preference tunability, we vary  $\lambda$  in increments of 0.1, capturing significant behavioral shifts. As a base model, we use *Mistral-7B-Instruct-v0.3* (Jiang et al., 2023) (licensed apache-2.0), training on NVIDIA A100 GPUs with an 80/20 train/test split, and 3% for validation. We run one epoch at a batch size of 4 and stop training when validation loss converges.

Apart from the special domain dataset, we also use the PKU-SafeRLHF dataset (licensed cc-by-nc-4.0) for safety and helpfulness alignment experiments (Ji et al., 2024).

## 6 Results and Discussion

### 6.1 Single Domain Preference Tuning

We use the AV derived by aligning the model to generate responses at an expert-level within a given domain. It facilitates model editing which introduces a tunable parameter, allowing the user to control the proficiency level of the generated responses in a continuum. Consequently, one alignment vector is established for each domain, enabling the model to navigate and produce output across varying spectra of proficiency. This, in turn, also addresses **RQ1**.

Table 1 shows that simply adding instructions for specific expertise (i.e., prompting) does not significantly improve preference accuracy, while nearly doubles inference cost. Notably, the base model achieves high expert-level accuracy even with prompts from a different LLM (Claude-3-Sonnet), though it performs poorly in generic (0.31) and avoidance (0.15) categories. For MEAV, adding the AV at different  $\lambda$  values shifts the model’s likelihood of generating expert responses: negative  $\lambda$  reduces expertise (with avoidance at  $\lambda = -1.2$ ), while in the medical domain,  $\lambda = -0.7$  yields generic behavior and  $\lambda = 0.5$  produces full expertise.

Figure 2 illustrates the tunable nature of the preference expertise spectrum across all three domains.

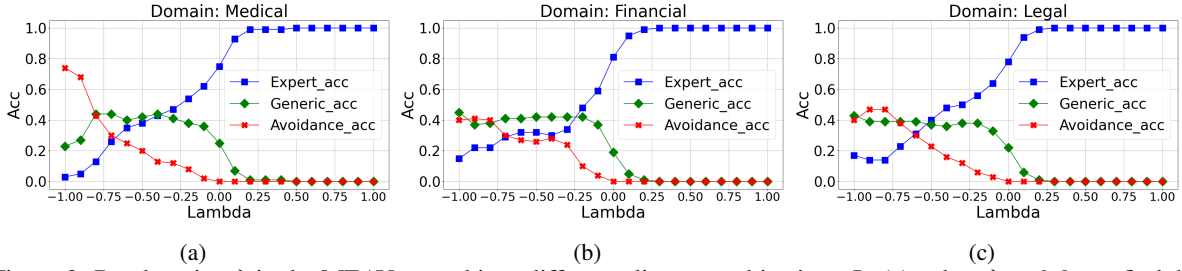


Figure 2: By changing  $\lambda$  in the MEAV, we achieve different alignment objectives. In (a), when  $\lambda > 0.3$ , we find the model aligning with expert answers to medical queries by preferring expert responses over the others. However, when  $\lambda < -0.8$ , we see the model prefers avoidance of responses. In between these points, we observe the model answering generically to medical queries. (b) and (c) demonstrates this behavior for financial and legal domains respectively.  $\lambda$  acts as a tunable knob to adjust model behavior.

Domain	Technique	Target behavior	Pref. Acc.			GPT-4 judged gen. acc		
			Exp	Gen	Avd	Exp	Gen	Avd
Medical	Default		.75	.25	0	.90	.05	.05
	Prompting	Exp	.78	.22	0	.90	.05	.05
		Gen	.69	.31	0	.50	.50	0
		Avd	.60	.25	.15	.15	.55	.30
Ours: MEAV	Exp (.5)	<b>.95</b>	0	.05	<b>1.0</b>	0	0	
	Gen (-.7)	.26	<b>.44</b>	.30	0	<b>.60</b>	.40	
	Avd (-1.2)	.03	.13	<b>.84</b>	.05	.20	<b>.75</b>	
Financial	Default		.81	.19	0	.85	.15	0
	Prompting	Exp	.84	.16	0	.95	.05	0
		Gen	.57	.43	0	.75	.25	0
		Avd	.35	.49	.16	.20	.60	.20
Ours: MEAV	Exp (.3)	<b>.85</b>	.15	0	<b>1.0</b>	0	0	
	Gen (-.4)	.30	<b>.42</b>	.28	.35	<b>.50</b>	.15	
	Avd (-1.4)	.07	.20	<b>.73</b>	0	.15	<b>.85</b>	
Legal	Default		.78	.22	0	.85	.15	0
	Prompting	Exp	.79	.21	0	1.0	0	0
		Gen	.59	.41	0	.65	.35	0
		Avd	.41	.30	.29	.15	.40	.45
Ours: MEAV	Exp (.3)	<b>1.0</b>	0	0	<b>1.0</b>	0	0	
	Gen (-.7)	.23	<b>.39</b>	.38	0	<b>.65</b>	.35	
	Avd (-1.4)	0	.20	<b>.80</b>	0	.05	<b>.95</b>	

Table 1: How MEAV performs to steer different domain expertise response level. The *Default* behavior indicates  $\lambda = 0$ , i.e., the model with no alignment. Tuning Lambda to different values with our MEAV approach leads to varying levels of proficiency responses. As such, we observe Exp, Gen, and Avd behavior just by aligning one model.

Notably, at  $\lambda = 0$ , the model predominantly generates expert responses in all domains. In the medical domain, the model reaches the higher end of the expertise spectrum when  $\lambda$  exceeds 0.3. Between  $\lambda = -0.4$  and  $\lambda = -0.8$ , the model exhibits varying degrees of generic behavior and beyond that, the model starts behaving with topic avoidance.

Next, we investigate if the gradual model editing method also impacts the performance in the other domains. Our findings indicate that the specialized behavior is indeed reflected across various domains, even when the AV is extracted for a specific domain. For instance, Table 2 demonstrates that the addition of a medical AV with  $\lambda = 0.5$  also enhances the model’s expertise in the financial do-

main. Similarly, we observed that with  $\lambda = -1.2$  the model exhibits avoidance behavior in both the legal and financial domains. This pattern is consistent when using other specialized domain vectors as well (see Appendix G).

**Effect on General Alignment** We also examine whether MEAV for controllable proficiency levels influences the general domain preference (i.e., ‘helpfulness’ and ‘safety’). Notably, we do not observe any regression in the safety domain; however, the model becomes increasingly helpful as  $\lambda$  increases. With the rise in  $\lambda$ , the model provides more detailed and specific guidance, which aligns with human preferences for helpfulness. Con-

Lambda	Fin pref. Acc			Leg pref. Acc			General Pref. Acc			
	Exp	Gen	Avd	Exp	Gen	Avd	Safety		Helpfulness	
							Safe	Unsafe	Helpful	Unhelpful
0	.81	.19	0	.78	.22	0	.58	.42	.60	.40
0.5	1.0	0	0	1.0	0	0	.58	.42	.66	.34
-0.7	.59	.40	.01	.58	.32	.10	.57	.43	.58	.42
-1.2	.03	.20	.77	.08	.18	.74	.57	.43	.49	.51

Table 2: Out of Domain (special and general) preference accuracy for Medical domain responses. Here, we gradually add the in-domain AV with the base model, and observe the performance for out-of-domain proficiency levels. We find that steering the proficiency levels in one domain also generalizes across other domains.

versely, decreasing  $\lambda$  causes the model to avoid answering, which is perceived as unhelpful. Notably, the range of change in general domain preference accuracy is  $\pm 11\%$  for helpfulness and  $\pm 1\%$  for safety, indicating that MEAV does not lead to significant regression in general domain performance.

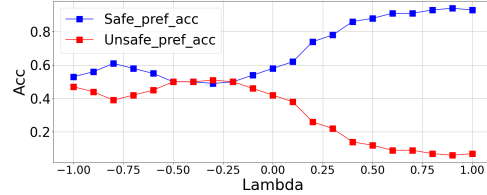


Figure 3: Controlling safety by MEAV

## 6.2 Multi Domain Preference Tuning

We observe distinct behaviors across different domains by adjusting specific configurations. Since, we have three proficiency levels, accuracy higher than 33% and the highest among the three levels can be considered as the “dominant” proficiency level. For example, as shown in Table 3, we find that an AV-based editing coefficient of -1, -1, and 0.6 for the Medical, Financial, and Legal domains, respectively, results in *avoidance* being the dominant behavior in the Medical and Financial domains, with accuracies of 0.46 and 0.42, respectively, and *expertise* being dominant in the Legal domain, with an accuracy of 0.78. Therefore, it indicates multi-level expertise across domains, and we address **RQ2** as well.

There are 27 possible domain-behavior combinations (three domains  $\times$  three spectrums), and a grid search reveals 22 where the desired behavior is dominant. Joint training achieves near-perfect accuracy but requires 27 separate trainings, nine times more than the three needed for single-domain DPO runs. Each training job takes about 72 hours on an A100 GPU, totaling 1,944 hours for all 27. By contrast, a grid search of 21 coefficient values per domain (9,261 evaluations at roughly 60 seconds each) takes about 155 hours, or 12 times faster. However, continuous multi-domain tunability remains challenging, as single-domain edits often over-generalize and compromise domain-specific precision.

However, this behavior can be effectively controlled by tuning the  $\lambda$  value. For example, in Figure 2a, for the medical expertise, the expert

spectrum spans from 0.3 to 1.0. When multidomain alignment is required, one can choose a lower lambda value to avoid overgeneralization while still achieving domain-specific attributes. Since the proficiency levels occur in a spectrum of  $\lambda$  rather than a fixed point, one can easily choose a suitable point based on their requirements.

## 6.3 Can AV be extensible for General Domain?

To explore the generalizability of MEAV across various domains, we focus on the safety alignment as a test case. We start by aligning our base model towards the *safe* dimension by obtaining the safety AV and gradually integrating it with the base model. For the safety alignment, we sample the examples where chosen response is labeled safe, and the rejected response is labeled unsafe (Ji et al., 2024). We compute the pref. acc in the same way described in 5.1, where  $R = \{safe, unsafe\}$ .

Figure 3 illustrates that the model exhibits mixed safety accuracy initially when  $\lambda = 0$  with a safety preference accuracy of 0.53 and an unsafe preference of 0.47. As  $\lambda$  increases, the model progressively aligns more with safety, achieving a safety preference accuracy of 0.93 at  $\lambda=1$ . However, when  $\lambda$  is adjusted negatively, the safety scores become inconsistent and mixed. Notably, even at large negative  $\lambda$  values, beyond -0.25, the model does not become fully “unsafe”.

In constructing the response proficiency levels, we intentionally maintain three distinct spectrums. In contrast, the PKU-SafeRLHF dataset does not

Baseline: Joint training			Ours: MEAV			editing coef
Med	Fin	Leg	Med	Fin	Leg	
Avd (100%)	Avd (99%)	Exp (98%)	Avd (46%)	Avd (42%)	Exp (78%)	[-1, -1, .6]
Avd (100%)	Exp (91%)	Exp (94%)	Avd (43%)	Exp (44%)	Exp (80%)	[-1, .8, .6]
Avd (100%)	Exp (90%)	Avd (90%)	Avd (57%)	Exp (56%)	Avd (36%)	[-.4, .4, -.8]
Exp (99%)	Avd (100%)	Exp (97%)	Exp (88%)	Avd (44%)	Exp (87%)	[.2, -.8, -.2]

Table 3: MEAV enables multidomain expertise through inference-time adjustment, avoiding the need to retrain separate models for each domain configuration. Unlike joint training, MEAV flexibly steers domain behavior using learned editing coefficients.

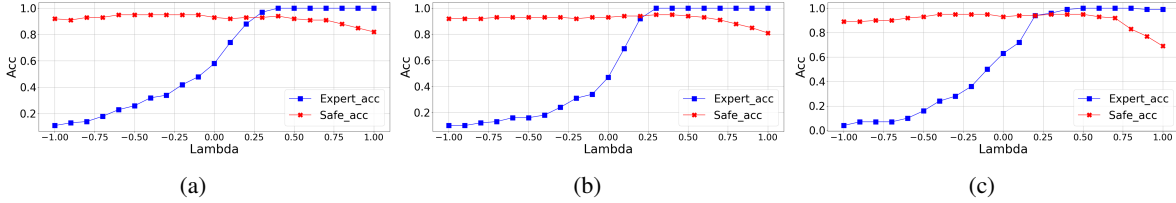


Figure 4: Visualizing the transferability of the MEAV process. We observe the effect of proficiency-level-encoded AV integration with a safety-aligned model in the (a) Medical (b) Financial (c) Legal domain proficiency control. For all domains within the range of -1 to 0.7, no safety regression happens, indicating the robustness of MEAV.

553 follow this structure, as it lacks any specific grada-  
554 tion in safety levels.

#### 555 6.4 Analyzing the Transferability of 556 Alignment Vector

557 Next, we explore whether AVs derived from a spe-  
558 cific alignment objective can be effectively applied  
559 to a pre-aligned model. As a case study, we select a  
560 safety-aligned version of the base model, to assess  
561 the transferability of these alignment vectors. Us-  
562 ing a similar approach to single-domain MEAV, we  
563 gradually integrate the AVs into our target model,  
564 which is safety-aligned.

565 Figure 4 presents the model’s performance as  $\lambda$   
566 is varied. Our findings indicate that when  $\lambda$  is ad-  
567 justed from -1 to +1, the model’s behavior related  
568 to safety, which is its primary control objective, re-  
569 mains relatively stable. For instance, in the medical  
570 domain (Figure 4(a)), varying  $\lambda$  within this range  
571 results in a minimal change in safety preference  
572 accuracy, with a difference of only 0.11 between  
573 the lowest and highest accuracy points. In contrast,  
574 the accuracy of medical expert response prefer-  
575 ences improves significantly, with an increase of  
576 0.81, which is over seven times greater than the  
577 change in safety preference accuracy. Hence, the  
578 AV obtained by our method is transferable to models  
579 aligned on other orthogonally aligned objectives as  
580 well, proving the transferability of MEAV.

#### 581 6.5 Generalizability Across Models

582 To assess the robustness and generalizability of  
583 MEAV, we extend our experiments beyond the  
584 Mistral-7B-Instruct model and evaluate its appli-  
585 cability to Qwen-2.5-3B-Instruct, a smaller-scale  
586 model with a distinct architectural profile (Bai et al.,  
587 2023). We found that by varying  $\lambda$  in the MEAV  
588 framework, a consistent transition across expert,  
589 generic, and avoidant response modes could be  
590 achieved. For example, in the medical domain,  
591  $\lambda$  value of 0.6, -0.4 and -1.40 can achieve Exp,  
592 Gen and Avd respectively (see more in appendix  
593 F). This suggests that our editing paradigm is not  
594 model-specific but broadly applicable across LLMs  
595 of different scales and architectures.

#### 596 7 Conclusion

597 We address inference-time preference alignment  
598 tunability through a novel model editing technique  
599 called MEAV. We build a synthetic dataset designed  
600 to represent three levels of response proficiency  
601 across three specialized domains. Our approach  
602 enables single-domain preference tunability at in-  
603 ference time without incurring additional costs or  
604 resource usage. This allows users to select differ-  
605 ent response proficiency levels without the need  
606 for extra training. Furthermore, our method offers  
607 tailored configurations for diverse multidomain be-  
608 haviors, significantly reducing both training time  
609 and resource demands.

## 610 **Limitations**

611 Our work has several limitations and areas for fu-  
612 ture exploration.

- 613 • We did not evaluate the correctness of the  
614 specialized domain responses. While the au-  
615 thors manually fact-checked a subset of the  
616 responses, and we also reported cross-LLM  
617 fact checks, we do not recommend using these  
618 synthetic LLM-generated responses without  
619 expert validation. Researchers found a 4.6%  
620 rate of hallucinations in Claude-generated re-  
621 sponse (Vectara, 2025). However, how the  
622 hallucinations might impact the special do-  
623 main responses, is left for future research.
- 624 • We used a basic approach (AV) for obtaining  
625 alignment vectors, which was simple and ef-  
626 fective for our use-case. However, whether  
627 the AVs are also capturing noise outside the  
628 preference dimension, is not explored in our  
629 work. To that end, more advanced techniques  
630 like parameter thresholding, zeroing, or SVD-  
631 based separation will be explored (Yadav et al.,  
632 2024; Gao et al., 2024) in our future work.
- 633 • Our method is currently applicable only to  
634 LLMs with the same architecture and parame-  
635 ter count. As new models with diverse archi-  
636 tectures and varying parameter sizes continue  
637 to emerge, this constraint may limit the gener-  
638 alizability of our approach. We aim to extend  
639 our methodology to support cross-architecture  
640 and cross-parameter adaptation in future.
- 641 • We relied on an extensive grid search for mul-  
642 tidomain alignment, which, while more ef-  
643 ficient than full retraining, remains compu-  
644 tationally intensive. A more optimized or  
645 strategic search approach could significantly  
646 reduce the parameter search space and further  
647 enhance efficiency.

## 648 **Ethical Implication and Broader Impact**

649 The introduction of MEAV offers a transformative  
650 approach to LLM alignment, enabling dynamic,  
651 inference-time preference adjustments while sig-  
652 nificantly reducing computational costs. This flex-  
653 ibility allows LLMs to be more adaptable across  
654 different specialty domains, such as medical, legal,  
655 and financial, without the need for retraining. How-  
656 ever, there are also some concerns with this, and  
657 we discuss this below:

- A model originally fine-tuned for safety-  
aligned behavior could be easily modified at  
inference time using adversarially crafted AVs  
to produce harmful, deceptive, or unsafe out-  
puts.
- The expert responses may encode cultural bias  
in all medical, legal, and financial domains.
- The ability to dynamically adjust model be-  
havior raises concerns about accountability,  
as users can shift LLM responses in ways that  
deviate from the ethical constraints originally  
intended.

## 670 **References**

- Anthropic. 2024. [Introducing the next generation of claude: The claude 3 family](#). Accessed: 2024-09-10.
- Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. 2024. A general theoretical paradigm to understand learning from human preferences. pages 4447–4455. PMLR.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. Scaling synthetic data creation with 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.
- Yi Dong, Zhilin Wang, Makesh Narsimhan Sreedhar, Xianchao Wu, and Oleksii Kuchaiev. 2023. Steerlm: Attribute conditioned sft as an (user-steerable) alternative to rlhf. *arXiv preprint arXiv:2310.05344*.
- Lei Gao, Yue Niu, Tingting Tang, Salman Avestimehr, and Murali Annavam. 2024. Ethos: Rectifying language models in orthogonal parameter space. *arXiv preprint arXiv:2403.08994*.
- Yiju Guo, Ganqu Cui, Lifan Yuan, Ning Ding, Jiexin Wang, Huimin Chen, Bowen Sun, Ruobing Xie, Jie Zhou, Yankai Lin, et al. 2024. Controllable preference optimization: Toward controllable multi-objective alignment. *arXiv preprint arXiv:2402.19085*.

709	James Y Huang, Sailik Sengupta, Daniele Bonadiman, Yi-an Lai, Arshit Gupta, Nikolaos Pappas, Saab Mansour, Katrin Kirchoff, and Dan Roth. 2024. Deal: Decoding-time alignment for large language models. <i>arXiv preprint arXiv:2402.06147</i> .	763
710		764
711		765
712		766
713		767
		768
714	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. <a href="#">Editing models with task arithmetic</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	769
715		770
716		771
717		772
718		
719	Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. 2023. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. <i>arXiv preprint arXiv:2310.11564</i> .	773
720		774
721		775
722		776
723		777
724		
725	Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2024. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. <i>Advances in Neural Information Processing Systems</i> , 36.	778
726		779
727		780
728		781
729		782
730		783
		784
731	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	785
732		786
733		787
734		788
735		789
736	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2024. Inference-time intervention: Eliciting truthful answers from a language model. <i>Advances in Neural Information Processing Systems</i> , 36.	790
737		791
738		792
739		793
740		794
		795
741	Lei Li, Yongfeng Zhang, and Li Chen. 2023. Prompt distillation for efficient llm-based recommendation. In <i>Proceedings of the 32nd ACM International Conference on Information and Knowledge Management</i> , pages 1348–1357.	796
742		797
743		798
744		799
745		800
746	Tianlin Liu, Shangmin Guo, Leonardo Bianco, Daniele Calandriello, Quentin Berthet, Felipe Linares, Jessica Hoffmann, Lucas Dixon, Michal Valko, and Mathieu Blondel. 2024. Decoding-time re-alignment of language models. <i>arXiv preprint arXiv:2402.02992</i> .	801
747		
748		802
749		803
750		
751		
752	Bertalan Meskó. 2023. Prompt engineering as an important emerging skill for medical professionals: tutorial. <i>Journal of medical Internet research</i> , 25:e50638.	804
753		805
754		806
		807
		808
755	OpenAI. 2025a. <a href="#">Gpt-4o system card</a> . Accessed: 2025-03-28.	809
756		810
		811
757	OpenAI. 2025b. <a href="#">o3-mini system card</a> . Accessed: 2025-03-28.	812
758		813
759	Jonas Oppenlaender, Rhema Linder, and Johanna Silvennoinen. 2023. Prompting ai art: An investigation into the creative skill of prompt engineering. <i>arXiv preprint arXiv:2303.13534</i> .	814
760		815
761		816
762		817
		818
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	
	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36.	
	Alexandre Rame, Guillaume Couairon, Corentin Dancette, Jean-Baptiste Gaya, Mustafa Shukor, Laure Soulier, and Matthieu Cord. 2023. Rewarded soups: towards pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards. <i>Advances in Neural Information Processing Systems</i> , 36:71095–71134.	
	Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2024. A systematic survey of prompt engineering in large language models: Techniques and applications. <i>arXiv preprint arXiv:2402.07927</i> .	
	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. <i>Advances in Neural Information Processing Systems</i> , 33:3008–3021.	
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	
	Vectara. 2025. <a href="#">Hallucination evaluation leaderboard</a> . Hugging Face Spaces. Accessed: 2025-02-15.	
	Pengyu Wang, Dong Zhang, Linyang Li, Chenkun Tan, Xinghao Wang, Ke Ren, Botian Jiang, and Xipeng Qiu. 2024. Inferaligner: Inference-time alignment for harmlessness through cross-model guidance. <i>arXiv preprint arXiv:2401.11206</i> .	
	Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. 2024. Ties-merging: Resolving interference when merging models. <i>Advances in Neural Information Processing Systems</i> , 36.	
	Rui Yang, Xiaoman Pan, Feng Luo, Shuang Qiu, Han Zhong, Dong Yu, and Jianshu Chen. 2024. Rewards-in-context: Multi-objective alignment of foundation models with dynamic preference adjustment. <i>arXiv preprint arXiv:2402.10207</i> .	

819 Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin  
820 Li. 2024. Language models are super mario: Absorb-  
821 ing abilities from homologous models as a free lunch.  
822 In *Forty-first International Conference on Machine*  
823 *Learning*.

824 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan  
825 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,  
826 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024.  
827 Judging llm-as-a-judge with mt-bench and chatbot  
828 arena. *Advances in Neural Information Processing*  
829 *Systems*, 36.

## 830 A Data generation and Annotation details

831 Table 4 shows the breakdown of the total amount  
832 of data collected.

Domain	Method		Total
	PH	CP	
Medical	5904	6096	13,000
Financial	6909	5465	12,374
Legal	5952	6915	12,867
Total curated data			38,241

Table 4: Full curated data amount. “PH” stands for PersonaHub and “CP” stands for CreatePersona

833 Table 5 shows the annotation accuracy for the  
834 human volunteers.

GT	Prediction		
	Exp	Gen	Avd
Exp	82.96%	17.04%	0%
Gen	17.04%	81.70%	1.26%
Avd	0%	1.26%	98.73%

Table 5: Average annotation accuracy for three annotators

## 835 B Synthetic Data Generation: How did 836 we arrive at the reported numbers of 837 generated data?

838 We evaluated the validity of persona-query pairs  
839 by manually reviewing a sample of 50 entries. Our  
840 analysis confirmed that Claude-3-sonnet reliably  
841 adhered to the instructions outlined in our prompt.  
842 To determine the dataset size, we initially generated  
843 15,000 queries across all domains, as described in  
844 Section 4.1.

845 During preprocessing, we identified and re-  
846 moved 1–3% of the queries and responses that were  
847 truncated due to timeout or quota limit issues. Ad-  
848 ditionally, queries in non-English languages were

849 excluded to ensure linguistic uniformity. We fur-  
850 ther filtered for completeness, retaining only those  
851 entries that contained a full set of persona-query  
852 responses across the three proficiency levels. Con-  
853 sequently, a small number of incomplete entries  
854 were discarded.

855 Next, we examined the dataset for query repe-  
856 tition. Although some personas were repetitive,  
857 we found no significant repetition in the queries  
858 themselves. Following this rigorous cleaning and  
859 filtering process, we finalized the dataset with the  
860 following distributions: 13,000 personas for the  
861 medical domain, 12,374 personas for the financial  
862 domain, and 12,867 personas for the legal domain.

## 863 C Persona and Query Generation 864 Prompts

865 The following prompts are used to generate Per-  
866 sonas and Queries:

### 867 C.1 CreatePersona

868 In the CreatePersona method, we generate a ran-  
869 dom root persona, and from that other persona are  
870 generated. Here is the prompt used:

871 **Prompt:** Based on the given persona, generate  
872 5 persona, that can be closely or remotely related  
873 to the given persona...generate a one-paragraph fi-  
874 nancial query in first person...

875 *Given Persona:* A patient coordinator with ex-  
876 cellent interpersonal skills and a knack for digital  
877 scheduling systems

#### 878 Claude-3-Sonnet generated query:

879 "persona1": "A recent college graduate with stu-  
880 dent loans and a entry-level job"

881 "query1": "Hello, I'm Alex, a 22-year-old who  
882 just graduated ...marketing position at a local firm,  
883 earning \$45,000 per year. However, I have accumu-  
884 lated \$32,000 in student loans ...financial situation  
885 and long-term objectives?"

886 "persona2": "A single mother juggling multiple  
887 part-time jobs and struggling to make ends meet"

888 "query2": "My name is Emily, and I'm a 32-  
889 year-old single mom working two part-time jobs  
890 to support my 6-year-old son...have accumulated  
891 over \$15,000 in credit card debt...increase my in-  
892 come or reduce expenses?"

### 893 C.2 PersonaHub

894 We generate queries directly from a given persona:

895 **Prompt:** Based on the persona described below,  
896 generate a one-paragraph medical query in first per-  
897 son, that the person fitting the persona can ask to an

898 online medical/health portal. Make sure the query  
899 is clear and very specific with nitty-gritty details  
900 like names, numbers etc, but brief. It should also in-  
901 clude relevant context, concerns, and other details  
902 to help the advisor or expert answer properly.

903 **Persona:** A retired coach known for their strate-  
904 gic approach to training and injury prevention

905 **Claude-3-Sonnet generated query:** As a re-  
906 tired coach ...I have a concerning issue that requires  
907 professional medical guidance. Over the past few  
908 weeks, I've been experiencing persistent discom-  
909 fort and stiffness in my left knee, particularly ...  
910 to regain full mobility and continue enjoying an  
911 active retirement without compromising my long-  
912 term joint health.

## 913 **D Experimental details**

- 914 • We used python3.11
- 915 • All experiments are done using pytorch 2.0
- 916 • Transformers version 4.43.0 was used

## 917 **E How Annotators were recruited**

918 We asked our authors (except for the first author)  
919 to volunteer as annotators. Each annotators were  
920 given detailed instructions with definitions of each  
921 response proficiency levels. The task was to select  
922 the responses for their respective levels.

### 923 **E.1 How instructions are provided to the** 924 **Annotators**

## 925 **Human Evaluation of Generated** 926 **Text**

927 We have provided you with the following special-  
928 ized domain query-response pairs. Your tasks are  
929 as follows:

930 **Task:** Rank the generated responses in the  
931 expertise level according to the level definitions.

932  
933 We generated three levels of responses: **Avoid-**  
934 **ance**, **Generic**, and **Expert**, for each query, by  
935 providing instructions based on the definitions  
936 below.

### 937 **Definitions:**

938 **(Level 1) Avoidance:** The response completely  
939 avoids providing any advice, information, or sug-  
940 gestions. It does not answer the question directly  
941 or indirectly. Instead, it politely informs to be  
942 unable to assist with the inquiries and recommends  
943 consulting an expert professional. Furthermore, the  
944

945 responder refrains from engaging in any discussion  
946 or providing any related resources or opinions  
947 regarding the issues specific to the domain in  
948 question.

949 **(Level 2) Generic:** The response represents  
950 general, non-specific information that can help  
951 guide the user but does not offer detailed advice. It  
952 offers general knowledge that is commonly known  
953 and easily accessible. The response is informative  
954 yet cautious. It emphasizes the importance of  
955 consulting an expert professional for accurate  
956 information and help and encourages users to seek  
957 advice from a qualified professional or institution  
958 for specific concerns.  
959

960  
961 **(Level 3) Expert:** The response is detailed, with  
962 expert-level advice and information. It thoroughly  
963 assesses the situations or context described and  
964 offers precise explanations and guidance tailored  
965 to the specific situation. The response reflects the  
966 depth and accuracy expected from an expert pro-  
967 fessional, and also the advice is not overly generic.  
968 Instead, it is comprehensive and nuanced, address-  
969 ing the user's specific circumstances. Finally, it  
970 offers clear, evidence-based recommendations and  
971 ensures the guidance is actionable and comprehen-  
972 sive.

973 **Instruction:** You will be given three responses  
974 for each query. You need to provide the ranking  
975 of each response separated by commas. For ex-  
976 ample, if you think Response 1 is Generic (level  
977 2), Response 2 is Expert (level 3), and Response  
978 3 is Avoidance (level 1), you should only answer:  
979 **2,3,1.**

980 You can also add a note if you want to notify us  
981 of something.

982 You will be provided with a spreadsheet with all  
983 these columns.

## 984 **F Generalizability across Models**

985 Our experiments suggest that the MEAV method  
986 of inference time alignment works for the smaller  
987 models as well. We experimented with Qwen-  
988 2.5-3B model, and found the similar performance  
989 across all domains. Check the Table 6 for more  
990 detail.

## 991 **G Out-of-domain performance**

992 Tab 7 and 8 reports the Out-of-domain performance  
993 for all three domains.

Lambda	Med pref. Acc			Fin pref. Acc			Leg Pref. Acc		
	Exp	Gen	Avd	Exp	Gen	Avd	Exp	Gen	Avd
[.6, -.4, -1.4]	.91	.09	.00	.74	.21	.05	.78	.22	.00
[.5, -.4, -1.6]	.35	.40	.25	.30	.50	.20	.13	.52	.35
[.5, -.3, -1.2]	.05	.27	.68	.04	.28	.68	.00	.23	.77

Table 6: MEAV performance for Qwen-2.5 3B-Instruct. Similar to the Mistral model, we observe that varying the Lambda value effectively steers the domain expertise level from Expert to Avoidant. The Lambda values shown to the left of each column correspond to the Medical, Financial, and Legal domains, respectively. These values were selected for demonstration purposes, and the resulting expertise behavior lies on a spectrum rather than being discrete.

Lambda	Med pref. acc			Leg pref. acc			Gen pref. acc			
	Exp	Gen	Avd	Exp	Gen	Avd	Safety		Helpfulness	
							Safe	Unsafe	Helpful	Unhelpful
0	.75	.25	0	.78	.22	0	.58	.42	.60	.40
.30	.97	.02	.01	.98	.02	0	.57	.43	.59	.41
-.40	.61	.37	.02	.57	.35	.08	.59	.41	.57	.43
-1.4	.18	.40	.42	.19	.52	.29	.55	.45	.51	.49

Table 7: Out of Domain (special and general) preference accuracy for Financial domain responses

Lambda	Med pref. acc			Fin pref. acc			Gen pref. acc			
	Exp	Gen	Avd	Exp	Gen	Avd	Safety		Helpfulness	
							Safe	Unsafe	Helpful	Unhelpful
0	.75	.25	0	.81	.19	0	.58	.42	.60	.40
.30	1.0	0	0	1.0	0	0	.53	.47	.59	.41
-.70	.30	.57	.13	.32	.56	.12	.56	.44	.53	.47
-1.4	.20	.58	.22	.13	.50	.37	.49	.51	.51	.49

Table 8: Out of Domain (special and general) preference accuracy for Legal domain responses