

Multi-View Attention Multiple-Instance Learning Enhanced by LLM Reasoning for Cognitive Distortion Detection

Jun Seo Kim^{1*} Hyemi Kim² Woo Joo Oh² Hongjin Cho² Hochul Lee² Hye Hyeon Kim^{3†}

¹Department of Computer Engineering, Gachon University

²AI Business Team, Korea Telecom Research

³Department of Biomedical Systems Informatics, Yonsei University

kma80kjs@gachon.ac.kr, {mika.kim, woojoo.oh, as.df, hochul.lee}@kt.com

hye_hyeon@yonsei.ac.kr

Abstract

Cognitive distortions have been closely linked to mental health disorders, yet their automatic detection remained challenging due to contextual ambiguity, co-occurrence, and semantic overlap. We proposed a novel framework that combines Large Language Models (LLMs) with Multiple-Instance Learning (MIL) architecture to enhance interpretability and expression-level reasoning. Each utterance was decomposed into Emotion, Logic, and Behavior (ELB) components, which were processed by LLMs to infer multiple distortion instances, each with a predicted type, expression, and model-assigned salience score. These instances were integrated via a Multi-View Gated Attention mechanism for final classification. Experiments on Korean (KoACD) and English (Therapist QA) datasets demonstrate that incorporating ELB and LLM-inferred salience scores improves classification performance, especially for distortions with high interpretive ambiguity. Our results suggested a psychologically grounded and generalizable approach for fine-grained reasoning in mental health NLP.

and illustrative examples of a subset of cognitive distortion types considered in this study. Influenced by internal factors such as emotions and beliefs, these distortions are expressed through language or automatic thoughts, reinforcing emotional distress and maladaptive behavior (Strohmeier et al., 2016). Cognitive distortions play a central role in the onset and maintenance of various mental illnesses, and early detection and intervention are considered key mechanisms of treatment (Morrison et al., 2015; Kaplan et al., 2017).

Recently, there have been active attempts to automatically detect cognitive distortions in mental health-related texts by utilizing the advanced language comprehension and reasoning capabilities of large language models (LLMs) (Chen et al., 2023; Qi et al., 2023). However, most existing studies treat utterances as single, unstructured inputs, returning predictions for the entire text without considering the internal psychological structure of the utterance. In particular, they overlook the fact that different cognitive distortions may arise from distinct aspects of an utterance—such as emotion, logic, or behavior—and that these components interact to shape distorted thinking. As a result, the interplay of such psychological factors often remains underrepresented, limiting interpretability and the granularity of model inference. Moreover, multiple cognitive distortions often occur together in a single utterance, and the semantic similarity between types can lead to differences in interpretation between experts or difficulties in establishing a gold standard (Suputra et al., 2023).

To overcome these limitations, we propose a new cognitive distortion detection framework inspired by the Multiple-Instance Learning (MIL) structure (Dietterich et al., 1997), in which an utterance is defined as a bag, and each of the multiple cognitive distortion expressions inferred by an LLM is considered as an instance to make a final decision. Each instance includes the predicted distortion

1 Introduction

Mental illness is a widespread global health concern. About half of the global population experiences mental illness in their lifetime, and one in eight is affected at any given time (McGrath et al., 2023; World Health Organization, 2022). Mental health conditions such as anxiety, depression, and emotional expression difficulties are significantly associated with cognitive distortions, indicating their role in both the formation and persistence of emotional distress (Mercan et al., 2023).

Cognitive distortion refers to systematic errors in thinking that occur when individuals perceive and interpret external information, leading to a negative conclusion that does not correspond to reality (Beck, 1979). Table 1 summarizes the definitions

Cognitive Distortion Type	Definition	Example
All-or-Nothing Thinking	Viewing situations in only two categories (e.g., perfect or failure) instead of on a spectrum.	“If I fail this test, I’m a total failure.”
Jumping to Conclusions	Predicting negative outcomes without evidence.	“She didn’t text back. She must be mad at me.”
Personalization	Blaming yourself for events outside your control or assuming excessive responsibility.	“My friend looks sad, maybe I did something wrong.”

Table 1: Definitions and examples of selected cognitive distortion types (adapted from Kim and Kim (2025)).

tion type, its associated sentence, and a salience score assigned by the LLM, which is incorporated into the final prediction through weighted aggregation within the MIL structure.

In addition, we designed this study to enable more precise and interpretable cognitive distortion reasoning by decomposing utterances into three psychologically grounded components—Emotion, Logic, and Behavior (ELB)—and inputting them into the LLM along with the original text. This approach moves beyond previous methods that rely solely on a single-text input, enabling more precise and interpretable prediction of complex, overlapping cognitive distortions in real-world utterances.

Our primary contributions are summarized as follows:

- We obtained high-quality labels for 10 cognitive distortion types through expert review by 10 psychologists.
- We structured each utterance into three psychologically grounded components (ELB) and incorporated this information into the LLM input to support more informed and context-aware inference.
- We proposed the first MIL-based framework that treats each LLM-inferred cognitive distortion as an instance and integrates both the predicted type and LLM-assigned salience score into a unified classification model.

2 Related Work

2.1 Cognitive Distortions Detection

Early studies treated utterances as single units and applied binary or multi-class classification using Linguistic Inquiry and Word Count features and models such as logistic regression or SVMs

(Simms et al., 2017; Shreevastava and Foltz, 2021). While effective in binary settings, these approaches struggled with label imbalance and semantic overlap in multi-class scenarios.

To capture co-occurring distortions, later work introduced multi-label classification (Ding et al., 2022; Shickel et al., 2020; Elsharawi and El Bolock, 2024), along with data augmentation and domain-adaptive language models. However, utterances were still processed holistically, without structural decomposition.

Some studies incorporated conversational context, modeling multi-turn interactions to improve continuity and prediction (Lybarger et al., 2022; Tauscher et al., 2023). Yet, these models also lacked expression-level inference and focused primarily on dialogue flow.

More recently, large LLMs have been applied to cognitive distortion detection. The Diagnosis of Thought (DoT) framework introduced a structured prompting approach to improve interpretability (Chen et al., 2023). Another study explored zero-shot and few-shot prompting for distortion classification without supervised training (Qi et al., 2023).

Despite recent progress, prior work did not model cognitive distortions at the expression level or incorporate the co-occurrence of multiple distortions within an utterance into prediction. To address this, we proposed a framework that decomposes utterances into ELB components, and infers distortions as instances within a Multiple-Instance Learning setup for more interpretable and fine-grained classification.

2.2 Multiple-Instance Learning in NLP

Multiple-Instance Learning (MIL) is a weakly supervised framework in which multiple instances are

grouped into a single bag and a prediction is made at the bag level. Originally proposed for drug activity prediction in bioinformatics (Dietterich et al., 1997), MIL was later applied to computer vision tasks such as natural scene classification (Maron and Ratan, 1998), demonstrating its flexibility in learning from partially labeled data.

In NLP, MIL has been applied to tasks such as document classification, sentiment analysis, and misinformation detection. Early approaches employed machine learning models such as mi-SVM, MILBoost, and other instance-level classifiers adapted to weakly supervised settings (Andrews et al., 2002; Zhang et al., 2008; Jorgensen et al., 2008).

Later work extended MIL with deep architectures to infer sentence-level sentiment from document-level labels. Approaches such as manifold regularization (Kotzias et al., 2014) and weighted instance modeling (Pappas and Popescu-Belis, 2014) improved both interpretability and prediction accuracy.

More recent studies have further enhanced MIL frameworks by integrating contextualized embeddings and attention mechanisms. For example, attention-based MIL has been applied to fake news detection with improvements in precision and interpretability (Karaoglan, 2024), and mutual-attention models have been used to address bag-instance mismatch in hate speech classification (Liu et al., 2022).

However, despite these advances, most MIL-based approaches in NLP define instances at the sentence or paragraph level, without leveraging finer-grained semantic representations. Incorporating LLM-inferred expressions and their types and salience scores into MIL has not yet been attempted. We address this gap by proposing a model that treats each LLM-generated unit as an instance and integrates its label and salience score into attention-based bag-level classification.

3 Dataset

KoACD The Korean Adolescent Cognitive Distortion (KoACD) dataset is derived from counseling texts collected on the NAVER Knowledge iN platform (Kim and Kim, 2025). For this study, we sampled 5,000 utterances (500 per distortion type) and conducted expert validation with 10 Korean psychologists holding master’s degrees and over five years of experience. Each utterance was reviewed

by a pair

Set	Utterances (%)
Train	3,608 (80%)
Validation	451 (10%)
Test	451 (10%)
Cognitive Distortion Type	Count (%)
Labeling	478 (10.6%)
Negative Filtering	470 (10.4%)
All-or-Nothing Thinking	464 (10.3%)
Emotional Reasoning	458 (10.2%)
Personalization	459 (10.2%)
Overgeneralization	452 (10.0%)
Discounting the Positive	451 (10.0%)
Jumping to Conclusions	431 (9.6%)
Magnification and Minimization	432 (9.6%)
Should Statements	415 (9.2%)
Total	4,510 (100.0%)

Table 2: Statistics of the KoACD Dataset

Set	Utterances (%)
Train	1,277 (80%)
Validation	159 (10%)
Test	161 (10%)
Cognitive Distortion Type	Count (%)
Mind Reading	239 (15.0%)
Overgeneralization	239 (15.0%)
Magnification	195 (12.2%)
Labeling	165 (10.3%)
Personalization	153 (9.6%)
Fortune-telling	143 (9.0%)
Emotional reasoning	134 (8.4%)
Mental filter	122 (7.6%)
Should statements	107 (6.7%)
All-or-nothing thinking	100 (6.3%)
Total	1,597 (100.0%)

Table 3: Statistics of the Therapist QA Dataset

of experts to cross-check the original labels, and after removing disagreements, 4,510 utterances with a single validated distortion label were retained. The data were split into training, validation, and test sets in an 8:1:1 ratio, and the label distribution is shown in Table 2.

Therapist QA Dataset This dataset consists of asynchronous patient–therapist Q&A logs from a publicly available Kaggle platform (Shreevastava and Foltz, 2021). We used a refined subset of 1,597 English-language utterances annotated by psychological experts with up to two of 10 cognitive distortion types, excluding utterances without distortions.

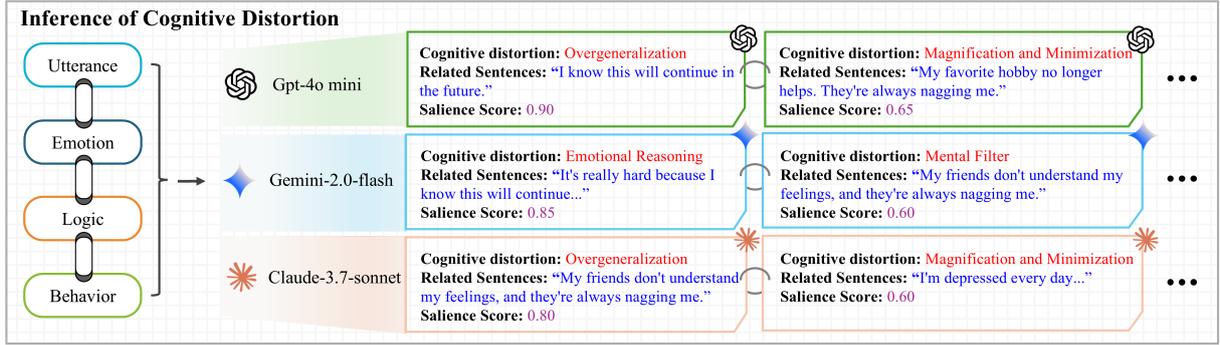


Figure 1: LLM-based Inference of Cognitive Distortion Instances from ELB-Structured Utterances

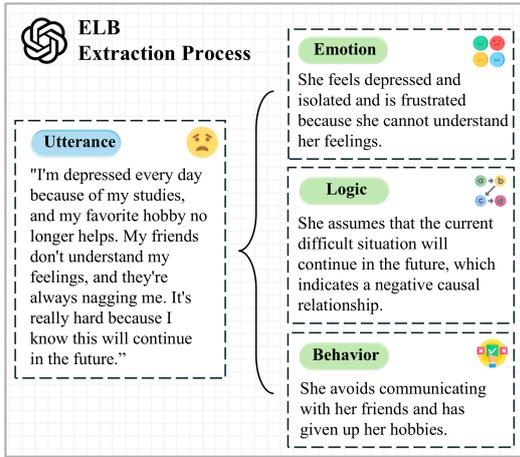


Figure 2: ELB-Based Psychological Decomposition of an Utterance

The dataset was used to benchmark cross-linguistic generalization with KoACD and was split into 8:1:1 for training, validation, and testing, as shown in Table 3.

4 Emotion–Logic–Behavior Extraction

To better capture the psychological context of each utterance, we decomposed it into three components—Emotion, Logic, and Behavior (ELB). This decomposition draws on the CBT cognitive triangle (Beck, 1979), with Logic used in place of Thought to emphasize its reasoning-oriented nature in cognitive distortions. This structured representation was designed to improve both the interpretability and granularity of cognitive distortion inference by making explicit the latent psychological dimensions, which are often entangled in unstructured text.

As illustrated in Figure 2, the ELB components were extracted using a zero-shot prompting strategy based on GPT-4 (OpenAI, 2023), which was guided to independently generate each of the three

elements for every utterance. These extracted components, combined with the original text, served as enriched input to the downstream LLM-based inference process, enabling more context-aware and psychologically grounded predictions.

The LLM hyperparameters are listed in Table 10 in Appendix B, and the prompt templates are provided in Table 16 in Appendix F.

5 LLM-Based Inference of Cognitive Distortion Instances

To infer multiple cognitive distortion instances per utterance, we employed three LLMs—OpenAI GPT-4o (OpenAI, 2023), Google Gemini 2.0 Flash (Google DeepMind, 2024), and Anthropic Claude 3.7 Sonnet (Anthropic, 2025). Each model independently processed the same utterance, using the pre-extracted ELB components as input to infer cognitive distortions. Each LLM produced a set of instances, consisting of a predicted distortion type, a corresponding text segment, and an LLM-assigned salience score, as illustrated in Figure 1.

Due to semantic ambiguity, a single sentence could map to multiple distortion types, either simultaneously or through varying interpretations. To evaluate the contribution of ELB information, we also performed inference using only the original utterances.

The LLM hyperparameters were listed in Table 10 in Appendix B, and the prompt templates were presented in Tables 17 and 18 in Appendix F.

6 Experiments

6.1 MIL Setup and Bag Construction

This study formulated cognitive distortion classification using a MIL framework, adopting an instance-based representation to effectively capture multiple distortions within a single utterance.

To this end, we aggregated cognitive distortion representations generated by three LLMs.

Instance Statistics	Value
Total Instances	52,201
Min Instances per Bag	5
Max Instances per Bag	20
Avg. Instances per Bag	11.57
Cognitive Distortion Type	# Instances (%)
Jumping to Conclusions	10,154 (19.5%)
Overgeneralization	9,942 (19.1%)
Negative Filtering	6,130 (11.7%)
Emotional Reasoning	6,095 (11.7%)
Personalization	5,531 (10.6%)
All-or-Nothing Thinking	4,739 (9.1%)
Labeling	3,580 (6.9%)
Magnification and Minimization	2,317 (4.4%)
Should Statements	2,190 (4.2%)
Discounting the Positive	1,523 (2.9%)
Total	52,201 (100.0%)

Table 4: Instance Distribution in KoACD (with ELB)

Instance Statistics	Value
Total Instances	13,967
Min Instances per Bag	4
Max Instances per Bag	24
Avg. Instances per Bag	8.75
Cognitive Distortion Type	# Instances (%)
Emotional reasoning	2,999 (21.5%)
Overgeneralization	2,271 (16.3%)
Fortune-telling	1,628 (11.7%)
Mind Reading	1,592 (11.4%)
Magnification	1,245 (8.9%)
All-or-nothing thinking	1,012 (7.2%)
Personalization	920 (6.6%)
Mental filter	899 (6.4%)
Should statements	783 (5.6%)
Labeling	618 (4.4%)
Total	13,967 (100.0%)

Table 5: Instance Distribution in Therapist QA (with ELB)

Each LLM independently inferred multiple cognitive distortion representations for a given utterance B , with each instance x_i defined by the predicted distortion type $type_i$, the associated expression text $text_i$, and a normalized confidence salience score \hat{s}_i assigned by the LLM. The entire utterance B is represented as a bag comprising a set of these instances. The overall construction is defined as shown in Equation 1:

$$B = \{x_i = (type_i, text_i, s_i)\}_{i=1}^N \quad (1)$$

Where N is the total number of instances generated by the model. Each instance’s salience score is normalized to reflect its relative importance within the model. The normalization process is defined as shown in Equation 2:

$$\hat{p}_i = \frac{s_i}{\sum_{j=1}^N s_j} \quad (2)$$

The resulting bag served as input to the MIL model, with instance-level salience score values used as weighting signals within the attention-based integration structure. Table 4 and Table 5 present instance-level summary statistics for the KoACD and Therapist QA datasets, respectively, under the ELB-based inference setting. These tables report the number of instances per cognitive distortion type, as well as the average and maximum number of instances per bag. For reference, statistics obtained without incorporating ELB information are provided in Table 8 and Table 9 in Appendix A.

6.2 Embedding Representation

Sentence Embedding Each bag consisted of a sentence embedding vector for the original utterance, which was combined with the instance embeddings. Utterances were embedded as 384-dimensional fixed-length vectors using the all-MiniLM-L12-v2 model (Wang et al., 2021). This embedding served two purposes: (1) it contributed to bag-level prediction when no instance captured the correct label, and (2) it preserved psychological and contextual information possibly missed at the instance level.

Instance Embedding Cognitive distortion instances generated from each utterance were also embedded using the same model. Each instance was formed by concatenating the predicted distortion type with the associated sentence and encoded into a 384-dimensional vector. For batch processing, bags were padded to the maximum number of instances observed in the dataset. The final input to the MIL model was a sequence combining the sentence embedding and all instance embeddings. Bag labels were one-hot encoded for loss computation only but were not included in the input embeddings.

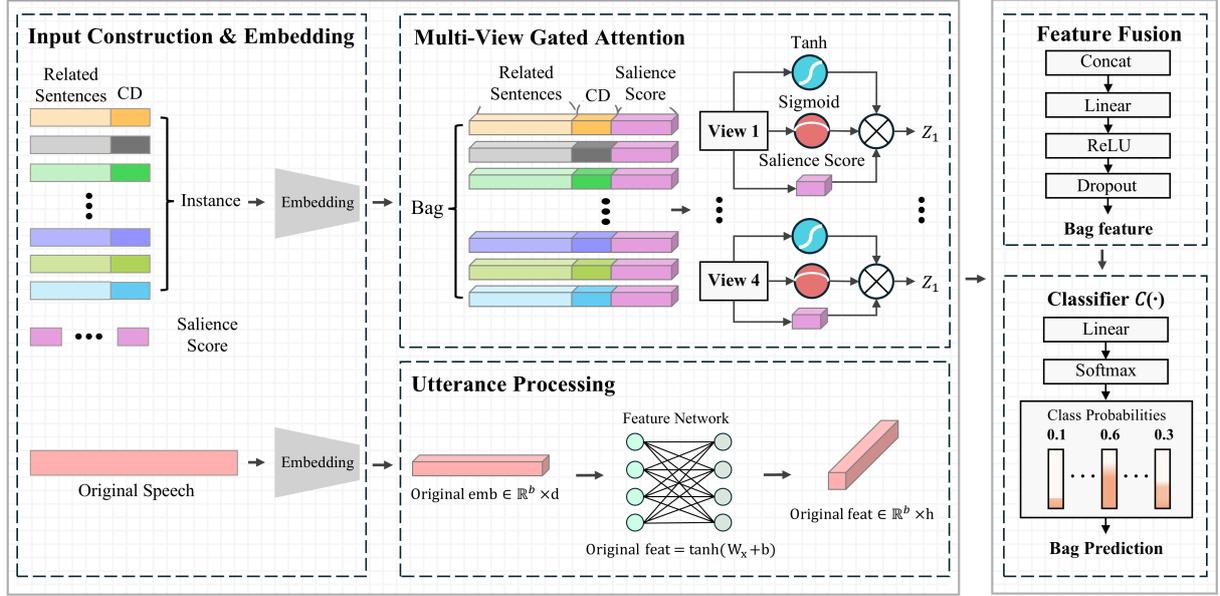


Figure 3: Multi-View MIL Architecture for Cognitive Distortion Classification

6.3 MIL with Multi-View Gated Attention Mechanism

We proposed a Multi-View Gated Attention model for cognitive distortion classification within the MIL framework. Unlike conventional MIL structures, our model integrated LLM-based instance representations and saliency scores into the attention mechanism. This architecture processed instances in parallel through multiple independent attention views and subsequently forms a unified representation for bag-level prediction, as illustrated in Figure 3.

Instance Weighting via Gated Attention. To compute the importance of each instance within the bag, we adopted a gated attention mechanism inspired by prior work on attention-based MIL (Ilse et al., 2018). The attention score for each instance i is defined as:

$$h_i = \sigma(W_g \cdot x_i) \cdot \tanh(W_f \cdot x_i) \cdot s_i \quad (3)$$

Here, x_i denotes the input vector of instance i , σ and \tanh represent the sigmoid and tanh activations, and W_g, W_f are learnable weight matrices. The sigmoid gate modulates the transformed feature vector from the tanh layer, and the result is scaled by the LLM-derived saliency score s_i . The final attention score h_i is computed independently for each view.

Multi-View Instance Modeling If attention is computed from a single view, it may fail to capture all

relevant instances, as only a subset may be attended to. To address this limitation, we adopt a multi-view attention structure inspired by Zhao et al. (2017). Each view independently computes instance-level attention scores using separate gate and feature networks. The final bag-level representation is obtained by averaging the outputs from all K views:

$$h_{\text{multi}} = \frac{1}{K} \sum_{k=1}^K h^{(k)} \quad (4)$$

Here, K denotes the total number of attention views, and $h^{(k)}$ is the aggregated instance representation from the k^{th} view.

Original Sentence Feature Transformation and Fusion. To incorporate global semantic information into the model, the original sentence embedding z is first transformed using a non-linear projection:

$$z' = \tanh(W_z \cdot z) \quad (5)$$

The aggregated instance-level representation h_{multi} from the multi-view attention is then concatenated with the transformed sentence embedding z' . This combined vector is passed through a linear projection followed by ReLU activation to produce the final bag-level representation v :

$$v = \text{ReLU}(W_c \cdot [h_{\text{multi}}, z']) \quad (6)$$

This final representation is used as input to a softmax classifier that predicts the type of cognitive distortion. Together, Equations 3–6 define the

core components of our model, which integrates instance-level attention with global sentence-level context.

6.4 Prediction and Evaluation

The final prediction was obtained by applying a softmax classifier over the fused bag-level representation. Model training was guided by the standard cross-entropy loss for multi-class classification. The learning rate was initialized at 0.0005 and decayed linearly by 0.00001 per epoch until reaching a minimum of 0.00001. Early stopping was applied if the validation loss did not improve for 10 consecutive epochs. Detailed model hyperparameters are provided in Table 11 in Appendix B. To evaluate generalization performance, all experiments were repeated with randomized train-validation-test splits. Results are reported as the mean \pm standard deviation of the weighted average F1 score across 10 runs.

7 Experimental Results and Analysis

7.1 Effect of ELB on Label Coverage

We defined a “missing” case as one in which none of the instances generated by the LLM included the ground-truth label of the utterance (bag). The missing rate served as a key metric for evaluating the model’s ability to capture expressions associated with the correct cognitive distortion.

As shown in Figure 4A, incorporating ELB information into the input reduced the missing rate across most distortion types in the KoACD dataset. Notable reductions were observed for *Magnification and Minimization* (from 35.42% to 27.08%) and *Emotional Reasoning* (from 11.35% to 5.90%). The overall average missing rate decreased from 10.89% to 8.93%, suggesting that explicitly structuring utterances into ELB components enabled the LLM to better identify label-relevant expressions.

Certain types, such as *Overgeneralization* and *Jumping to Conclusions*, maintained low missing rates regardless of ELB usage. In contrast, types like *Labeling* and *Discounting the Positive*, which are characterized by higher expressive variability and semantic ambiguity, continued to exhibit relatively higher omission rates.

Similar trends were observed in Figure 4B for the Therapist QA dataset. ELB usage led to reductions in missing rates for all distortion types, with substantial improvements in *Labeling* (from 38.18% to 30.30%), *Personalization* (from 32.68%

to 17.65%), and *Mental Filter* (from 43.44% to 32.79%). The largest reduction was observed for *Emotional Reasoning*, where the missing rate dropped from 17.91% to 2.99%.

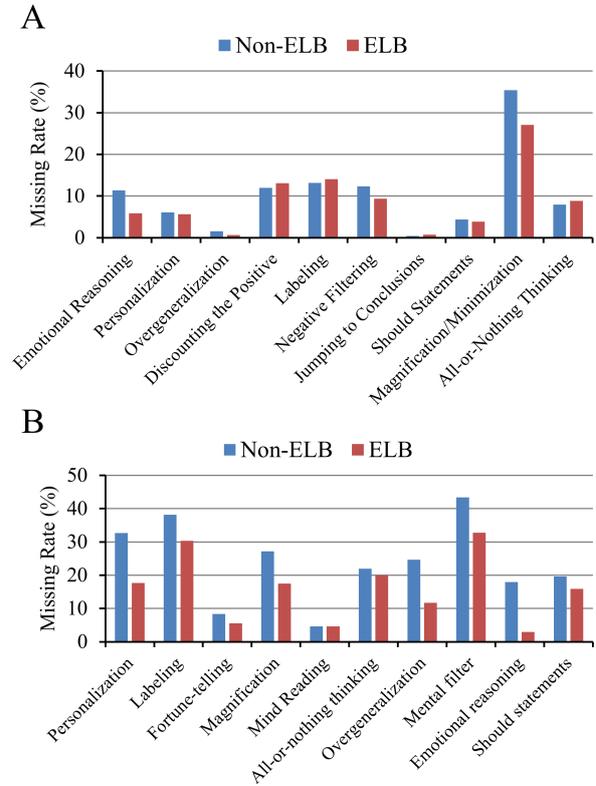


Figure 4: Comparative Missing Rates of Cognitive Distortion Instances With and Without ELB (A: KoACD, B: Therapist QA)

7.2 Input Configuration Comparison

We evaluated four input configurations with or without ELB components and LLM-inferred salience scores to analyze their impact on cognitive distortion classification performance. Four input conditions were set for this purpose:

1. **Baseline:** Neither ELB nor salience scores included,
2. **ELB only:** ELB components included, salience scores excluded,
3. **Salience only:** Salience scores included, ELB components excluded,
4. **ELB + Salience:** Both ELB components and salience scores included.

The results are summarized in Table 6 for both the KoACD and Therapist QA datasets. In the KoACD experiment, the configuration using both

Methods	KoACD		Therapist QA	
	Val F1	Test F1	Val F1	Test F1
Baseline	0.504 ± 0.019	0.473 ± 0.015	0.410 ± 0.038	0.340 ± 0.037
ELB	0.519 ± 0.016	0.483 ± 0.017	0.438 ± 0.028	0.378 ± 0.036
Saliency	0.518 ± 0.015	0.486 ± 0.014	0.428 ± 0.036	0.360 ± 0.035
ELB + Saliency	0.529 ± 0.018	0.505 ± 0.014	0.460 ± 0.029	0.394 ± 0.034

Table 6: Effect of ELB and Saliency Score Weighting on Classification Performance: All values are 10-run weighted F1 averages; best results are in bold.

ELB and LLM-derived saliency scores achieved the highest performance, recording a validation F1-score of 0.529 and a test F1-score of 0.505. Notably, the ELB-only configuration outperformed

Cognitive Distortion Type	KoACD (F1)	Therapist QA (F1)
Discounting the Positive	0.808 ± 0.036	-
Should Statements Labeling	0.852 ± 0.039	0.460 ± 0.084
Mind Reading	0.607 ± 0.052	0.469 ± 0.102
Personalization	-	0.608 ± 0.061
Overgeneralization	0.591 ± 0.036	0.409 ± 0.069
Jumping to Conclusions	0.440 ± 0.057	0.415 ± 0.077
Fortune-telling	0.427 ± 0.048	-
All-or-Nothing Thinking	-	0.437 ± 0.089
Magnification and Minimization (Magnification in Therapist QA)	0.456 ± 0.077	0.164 ± 0.097
Emotional Reasoning	0.390 ± 0.065	0.361 ± 0.120
Negative Filtering (Mental filter)	0.297 ± 0.052	0.164 ± 0.122
	0.217 ± 0.047	0.214 ± 0.119

Table 7: Per-Type F1 Scores on KoACD and Therapist QA: Top three scores for each dataset are shown in bold.

the saliency-only variant, suggesting that psychologically grounded structuring provides greater benefits than relying solely on LLM-inferred saliency scores.

A similar trend was observed in the Therapist QA dataset, where the combination of ELB and saliency scores again yielded the best performance. Compared to the DoT-based GPT-4 model with a test F1-score of 0.346 (Chen et al., 2023), our model achieved a higher score of 0.394, highlighting the effectiveness of incorporating structured psychological components. Across both datasets, configurations incorporating ELB consistently outperformed those relying solely on LLM-derived saliency scores, suggesting that structuring utter-

ances into psychological components enhances classification effectiveness.

7.3 Type-wise Performance Analysis

We evaluated classification performance by distortion type under the ELB + Saliency configuration, with the results for both KoACD and Therapist QA datasets summarized in Table 7.

Across both datasets, *Should Statements* consistently achieved the highest F1-scores. However, a substantial performance gap was observed for this type, with KoACD outperforming Therapist QA. Expert analysis attributes this discrepancy to linguistic style differences: KoACD contains concise, emotionally direct expressions, while Therapist QA includes longer narratives where obligation-related cues are embedded in complex contexts.

In contrast, distortion types involving emotionally ambiguous or abstract reasoning—such as *Emotional Reasoning* and *Magnification and Minimization*—showed lower F1-scores and greater variability, particularly in the English dataset. Further discussion and representative utterances can be found in Table 15 in Appendix E.

8 Conclusion and Future Work

This study proposed a novel framework for cognitive distortion detection by representing each utterance as a bag of LLM-generated instances, where each instance included a predicted distortion type, its associated expression, and a model-assigned saliency score. These components were integrated into a MIL architecture, allowing the model to attend to instances based on both semantic relevance and LLM-derived saliency scores.

By structuring utterances into ELB components and integrating them with LLM-inferred saliency scores, the framework enabled finer-grained inference and broader label coverage. These results highlight the potential of the proposed framework as a practical NLP tool for early detection and analysis of cognitive distortion.

Future work should focus on reducing omission rates and improving performance for distortion types characterized by low predictive accuracy or interpretive ambiguity. Notably, distortion types with higher instance frequency did not always correspond to better classification outcomes, suggesting the importance of incorporating quality-aware modeling strategies.

9 Limitations

In this study, incorporating ELB information into the input configuration significantly reduced the omission rate of correct labels for cognitive distortion expressions. However, certain limitations remain. True labels were not always captured at the instance level, particularly for distortion types that involve subtle emotional cues or subjectively phrased language. This may reflect both the inherent representational ambiguity of natural language and the challenges LLMs face in consistently detecting psychologically meaningful patterns across diverse linguistic expressions.

Moreover, although each cognitive distortion type was equally represented at the utterance level (500 samples per type), the number of inferred instances varied significantly across types. As a result, the MIL model disproportionately focused on distortion types with more abundant or salient instances during attention computation, potentially overlooking valid signals from types with sparser or less prominent instance patterns. These findings underscore the need for more balanced instance generation and attention regulation mechanisms in future work.

Furthermore, although the proposed MIL-based structure provides indirect interpretability through instance-level attention, it still lacks explicit and quantitatively grounded explanations for the model’s final decisions. Given that cognitive distortions are closely tied to clinical judgments, transparent psychologically consistent explanations of model predictions are essential for ensuring the model’s clinical applicability. From this perspective, future research should explore additional explanation generation structures or the integration of quantitative evidence frameworks grounded in psychological theory to further enhance interpretability.

10 Ethical Considerations

All data used in this study were obtained from publicly accessible platforms and were anonymized prior to publication. NAVER Knowledge and Kaggle public Q&A data were provided without personal identifiers (e.g., names, contact information, or account details), and no additional identifying information was collected or processed. Therefore, the data adhered to ethical standards of anonymity without requiring further de-identification.

To ensure the validity of the cognitive bias labels, 10 licensed psychological experts independently reviewed and cross-validated 5,000 utterances from the KoACD dataset. Each expert evaluated the appropriateness of the assigned cognitive distortion type based on predefined criteria. All utterances were fully anonymized to prevent any personal identification, and no personal data were collected from the annotators.

All annotators voluntarily participated in the study after being informed of the nature and purpose of the task. They were fairly compensated for their time and expertise, receiving 300,000 KRW for annotating 1,000 utterances—an amount calculated to reflect professional hourly rates based on estimated task duration. The annotation process was conducted under strict anonymity. This study adheres to the principles of responsible research involving human expert participants.

It is important to note that the cognitive distortion detection model proposed in this study is not intended as a substitute for a clinical diagnosis tool. Its use without oversight by qualified mental health professionals is not recommended. As LLM-based models may yield inaccurate or biased interpretations of emotionally sensitive language, real-world applications must involve ethical review frameworks and expert supervision to ensure safe and responsible deployment.

References

- Stuart Andrews, Ioannis Tsochantaridis, and Thomas Hofmann. 2002. Support vector machines for multiple-instance learning. In *Advances in Neural Information Processing Systems*, volume 15.
- Anthropic. 2025. [Claude 3.7 sonnet](#).
- Aaron T. Beck. 1979. *Cognitive therapy and the emotional disorders*. Penguin.
- Zhiyu Chen, Yujie Lu, and William Yang Wang. 2023. [Empowering psychotherapy with large language](#)

- models: Cognitive distortion detection through diagnosis of thought prompting. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4295–4304, Singapore. Association for Computational Linguistics.
- Thomas G. Dietterich, Richard H. Lathrop, and Tomás Lozano-Pérez. 1997. Solving the multiple instance problem with axis-parallel rectangles. *Artificial Intelligence*, 89(1-2):31–71.
- Xiruo Ding, Kevin Lybarger, Justin Tauscher, and Trevor Cohen. 2022. Improving classification of infrequent cognitive distortions: Domain-specific model vs. data augmentation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Student Research Workshop*, pages 68–75, Seattle, USA. Association for Computational Linguistics.
- Nada Elsharawi and Alia El Bolock. 2024. C-journal: A journaling application for detecting and classifying cognitive distortions using deep-learning based on a crowd-sourced dataset. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 3224–3234, Torino, Italia. ELRA and ICCL.
- Google DeepMind. 2024. Gemini 2.0 flash.
- Maximilian Ilse, Jakub Tomczak, and Max Welling. 2018. Attention-based deep multiple instance learning. In *Proceedings of the 35th International Conference on Machine Learning (ICML 2018)*, pages 2127–2136. PMLR.
- Zach Jorgensen, Yan Zhou, and Meador Inge. 2008. A multiple instance learning strategy for combating good word attacks on spam filters. *Journal of Machine Learning Research*, 9(6).
- S.C. Kaplan, A.S. Morrison, P.R. Goldin, R.T. Olino, and J. Gross. 2017. The cognitive distortions questionnaire (cd-quest): Validation in a sample of adults with social anxiety disorder. *Cognitive Therapy and Research*, 41:576–587.
- Kürsat Mustafa Karaoğlan. 2024. Novel approaches for fake news detection based on attention-based deep multiple-instance learning using contextualized neural language models. *Neurocomputing*, 602:128263.
- Jun Seo Kim and Hye Hyeon Kim. 2025. Koacd: The first korean adolescent dataset for cognitive distortion analysis. *Preprint*, arXiv:2505.00367.
- Dimitrios Kotzias, Misha Denil, Phil Blunsom, and Nando de Freitas. 2014. Deep multi-instance transfer learning. *Preprint*, arXiv:1411.3128.
- Jiexi Liu, Dehan Kong, Longtao Huang, Dinghui Mao, and Hui Xue. 2022. Multiple instance learning for offensive language detection. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7387–7396, Abu Dhabi, UAE. Association for Computational Linguistics.
- Kevin Lybarger, Justin Tauscher, Xiruo Ding, Dror Ben-Zeev, and Trevor Cohen. 2022. Identifying distorted thinking in patient-therapist text message exchanges by leveraging dynamic multi-turn context. In *Proceedings of the Eighth Workshop on Computational Linguistics and Clinical Psychology*, pages 126–136, Seattle, USA. Association for Computational Linguistics.
- Oded Maron and Aparna Lakshmi Ratan. 1998. Multiple-instance learning for natural scene classification. In *Proceedings of the 15th International Conference on Machine Learning (ICML 1998)*, pages 341–349.
- John J. McGrath, Ali Al-Hamzawi, Jordi Alonso, Yamin Altwaijri, Laura H. Andrade, Evelyn J. Bromet, Ronny Bruffaerts, José Miguel Caldas-de Almeida, Stephanie Chardoul, Wai Tat Chiu, Louisa Degenhardt, Olga V. Demler, Finola Ferry, Oye Gureje, Josep Maria Haro, Elie G. Karam, Georges Karam, Salma M. Khaled, Viviane Kovess-Masfety, and 40 others. 2023. Age of onset and cumulative risk of mental disorders: a cross-national analysis of population surveys from 29 countries. *The Lancet Psychiatry*, 10(9):668–681.
- Nihan Mercan, Merve Bulut, and Çiğdem Yüksel. 2023. Investigation of the relatedness of cognitive distortions with emotional expression, anxiety, and depression. *Current Psychology*, 42:2176–2185.
- Amanda S. Morrison, Carrie M. Potter, Matthew M. Carper, Dina G. Kinner, Dane Jensen, Laura Bruce, Judy Wong, Irismar Reis de Oliveira, Donna M. Sudak, and Richard G. Heimberg. 2015. The cognitive distortions questionnaire (cd-quest): Psychometric properties and exploratory factor analysis. *International Journal of Cognitive Therapy*, 8(4):287–305.
- OpenAI. 2023. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Nikolaos Pappas and Andrei Popescu-Belis. 2014. Explaining the stars: Weighted multiple-instance learning for aspect-based sentiment analysis. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 455–466, Doha, Qatar. Association for Computational Linguistics.
- Hongzhi Qi, Qing Zhao, Jianqiang Li, Changwei Song, Wei Zhai, Dan Luo, Shuo Liu, Yi Jing Yu, Fan Wang, Huijing Zou, Bing Xiang Yang, and Guanghui Fu. 2023. Supervised learning and large language model benchmarks on mental health datasets: Cognitive distortions and suicidal risks in chinese social media. *Preprint*, version 1.
- Benjamin Shickel, Scott Siegel, Martin Heesacker, Sherry Benton, and Parisa Rashidi. 2020. Automatic detection and classification of cognitive distortions in mental health text. In *Proceedings of the 2020 IEEE*

20th International Conference on Bioinformatics and Bioengineering (BIBE), pages 275–280. IEEE.

Sagarika Shreevastava and Peter W. Foltz. 2021. [Detecting cognitive distortions from patient-therapist interactions](#). In *Proceedings of the Seventh Workshop on Computational Linguistics and Clinical Psychology: Improving Access*, pages 151–158. Association for Computational Linguistics.

Taetem Simms, Christopher Ramstedt, Marc Rich, Matthew Richards, Thomas Martinez, and Christophe Giraud-Carrier. 2017. [Detecting cognitive distortions through machine learning text analytics](#). In *Proceedings of the 2017 IEEE International Conference on Healthcare Informatics (ICHI)*, pages 508–512. IEEE.

Craig W. Strohmeier, Brad Rosenfield, Robert A. DiTomasso, and J. Russell Ramsay. 2016. [Assessment of the relationship between self-reported cognitive distortions and adult adhd, anxiety, depression, and hopelessness](#). *Psychiatry Research*, 238:153–158.

I. Putu Gede Harsa Suputra, Linawati, Ni Putu Sastra, Gede Sukadarmika, Ni Kadek Anindya Saraswati, Diah Purwitasari Er, and I Made Agus Setiawan. 2023. [Detection and classification of cognitive distortions: A literature review](#). In *Proceedings of the 2023 International Conference on Smart-Green Technology in Electrical and Information Systems (ICS-GTEIS)*, pages 166–171. IEEE.

Justin S. Tauscher, Kevin Lybarger, Xiruo Ding, Ayesha Chander, William J. Hudenko, Trevor Cohen, and Dror Ben-Zeev. 2023. [Automated detection of cognitive distortions in text exchanges between clinicians and people with serious mental illness](#). *Psychiatric Services*, 74(4):407–410.

Wenhui Wang, Hangbo Bao, Shaohan Huang, Li Dong, and Furu Wei. 2021. [Minilmv2: Multi-head self-attention relation distillation for compressing pre-trained transformers](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2140–2151. Association for Computational Linguistics.

World Health Organization. 2022. World mental health report: Transforming mental health for all.

Qi Zhang, Sally A. Goldman, Wei Yu, and Jason E. Fritts. 2008. [Em-dd: An improved multiple-instance learning algorithm](#). In *Proceedings of the 19th IEEE International Conference on Pattern Recognition*, pages 1–4. IEEE.

Zhi Zhao, Yang Zhang, and Xindong Wu. 2017. Multiple instance learning via deep kernel embedding. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI 2017)*, pages 3442–3448. AAAI Press.

A Instance Statistics Without ELB Input

For completeness, we provide supplementary statistics summarizing the instance distribution for both datasets when ELB components are not included in the LLM input, as shown in Table 8 (KoACD) and Table 9 (Therapist QA). These tables offer a point of reference for understanding how instance generation differs in the absence of structured psychological input.

Instance Statistics	Value
Total Instances	49,229
Min Instances per Bag	4
Max Instances per Bag	19
Avg. Instances per Bag	10.92
Cognitive Distortion Type	# Instances (%)
Jumping to Conclusions	10,106 (20.6%)
Overgeneralization	8,809 (17.9%)
Negative Filtering	5,606 (11.4%)
Personalization	5,436 (11.0%)
All-or-Nothing Thinking	4,908 (10.0%)
Emotional Reasoning	4,728 (9.6%)
Labeling	3,711 (7.5%)
Should Statements	2,211 (4.5%)
Magnification and Minimization	2,169 (4.4%)
Discounting the Positive	1,545 (3.1%)
Total	49,229 (100.0%)

Table 8: Instance Distribution in KoACD (without ELB)

Instance Statistics	Value
Total Instances	10,788
Min Instances per Bag	2
Max Instances per Bag	21
Avg. Instances per Bag	6.76
Cognitive Distortion Type	# Instances (%)
Emotional Reasoning	2,024 (18.8%)
Mind Reading	1,528 (14.2%)
Fortune-telling	1,476 (13.7%)
Overgeneralization	1,297 (12.0%)
All-or-Nothing Thinking	1,015 (9.4%)
Magnification	938 (8.7%)
Mental Filter	696 (6.5%)
Personalization	693 (6.4%)
Should Statements	617 (5.6%)
Labeling	504 (4.7%)
Total	10,788 (100.0%)

Table 9: Instance Distribution in Therapist QA (without ELB)

B Experimental Settings and Model Hyperparameters

All experiments were conducted on a local machine equipped with an AMD Ryzen 9 7900X 12-Core Processor (4.70 GHz) without GPU acceleration. The implementation was based on PyTorch 2.6.0.

The hyperparameters used for all LLMs during both ELB extraction and cognitive distortion inference are summarized in Table 10. These values were applied uniformly across GPT-4o, Gemini 2.0 Flash, and Claude 3.7 Sonnet.

The model hyperparameters for the Multi-View Gated Attention MIL architecture are provided in Table 11. These settings were determined through a combination of grid search and empirical tuning based on validation performance and were consistently used across all experimental configurations.

Hyperparameter	Value
Max Tokens	512
Temperature	0.7
Top-p	1.0
Frequency Penalty	0.0
Presence Penalty	0.0

Table 10: LLM Hyperparameters for ELB Extraction and Cognitive Distortion Inference (applied to GPT-4o, Gemini 2.0 Flash, and Claude 3.7 Sonnet).

Hyperparameter	Value
Learning Rate	0.0005 (decayed 0.00001/epoch, min 0.00001)
Batch Size	32
Embedding Dimension	384
Attention Views (K)	4
Hidden Dimension	128
Dropout Rate	0.5
Optimizer	Adam
Early Stopping Patience	10 epochs

Table 11: Model Hyperparameters for the Multi-View Gated Attention MIL Architecture.

C Cognitive Distortion Types

This appendix summarizes the definitions and examples of the ten cognitive distortion types used in each dataset. Although both KoACD and Therapist QA adopt ten-type schemes, the specific categories partially overlap and partially diverge, reflecting linguistic and cultural differences. These definitions are presented in Table 12 for KoACD, taken

Cognitive Distortion Type	Definition	Examples
All-or-Nothing Thinking	Viewing situations in only two categories (e.g., perfect or failure) instead of on a spectrum.	“If I fail this test, I’m a total failure.”
Overgeneralization	Drawing broad conclusions from a single event or limited evidence.	“My one friend ignored me, so everyone else will hate me too.”
Mental Filtering	Focusing only on the negative aspects of a situation while ignoring the positive.	“I only remember my mistake though I got compliments on my presentation.”
Discounting the Positive	Rejecting positive experiences or compliments by insisting they don’t count.	“People told me I did well, but I was just being polite.”
Jumping to Conclusions	Predicting negative outcomes without evidence.	“She didn’t text back. She must be mad at me.”
Magnification and Minimization	Exaggerating negative or risky aspects while minimizing positive or positive aspects.	“One little mistake at work means I’m incompetent.”
Emotional Reasoning	Believing something must be true because you feel it strongly.	“I feel worthless, so I must be worthless.”
“Should” Statements	Holding rigid rules about how you or others should behave, leading to guilt or frustration.	“I should always be productive; otherwise, I’m lazy.”
Labeling	Assigning negative labels to yourself or others based on one event.	“I made a mistake, so I’m a total failure.”
Personalization	Blaming yourself for events outside your control or assuming excessive responsibility.	“My friend looks sad, maybe I did something wrong.”

Table 12: Cognitive Distortion Types in KoACD (adapted from [Kim and Kim \(2025\)](#))

from [Kim and Kim \(2025\)](#), and in Table 13 for Therapist QA, taken from [Shreevastava and Foltz \(2021\)](#).

Cognitive Distortion Type	Definition	Examples
Personalization	Personalizing or taking up the blame for a situation, that in reality involved many factors and was out of the person's control.	"My son is pretty quiet today. I wonder what I did to upset him."
Mind Reading	Suspecting what others are thinking or what are the motivations behind their actions.	"My house was dirty when my friends came over, they must think I'm a slob!"
Overgeneralization	Major conclusions are drawn based on limited information.	"Last time I was in the pool I almost drowned, I am a terrible swimmer and should not go into the water again."
All-or-Nothing Thinking	Looking at a situation as either black or white or thinking that there are only two possible outcomes to a situation.	"If I cannot get my Ph.D., then I am a total failure."
Emotional Reasoning	Giving someone or something a label without finding out more about it/them.	"Even though Steve is here at work late every day, I know I work harder than anyone else at my job."
Labeling	Emphasizing the negative or playing down the positive of a situation.	"My daughter would never do anything I disapproved of."
Magnification	Believing something must be true because you feel it strongly.	"My professor said he made some corrections on my paper, so I know I'll probably fail the class."
Mental Filter	Placing all one's attention on, or seeing only, the negatives of a situation.	"My husband says he wishes I was better at housekeeping, so I must be a lousy wife."
Should Statements	Should statements appear as a list of ironclad rules about how a person should behave, this could be about the speaker themselves or others. It is NOT necessary that the word "should" or its synonyms (ought to, must, etc.) be present in the statements containing this distortion.	"I should get all A's to be a good student."
Fortune-telling	Expecting things to happen a certain way, or assuming that things will go badly. Counterintuitively, this distortion does not always have future tense.	"I was afraid of job interviews, so I decided to start my own thing."

Table 13: Cognitive Distortion Types in Therapist QA (adapted from [Shreevastava and Foltz \(2021\)](#))

D Examples of Cognitive Distortion Instances Inferred by LLMs

We present an example of a single utterance (bag) from the KoACD dataset and the set of cognitive distortion instances inferred by three LLMs—GPT-4o, Gemini 2.0 Flash, and Claude 3.7 Sonnet. Each instance includes a predicted distortion type, the relevant text, and an LLM-assigned salience score, illustrating how a bag is composed of multiple instances in our experimental setting, as shown in Table 14.

E Cross-Dataset Variation in “Should Statements”

Despite sharing the same distortion label, “Should Statements” exhibited a substantial performance gap between the KoACD and Therapist QA datasets. This discrepancy can be attributed to the structural and contextual characteristics of each dataset.

The KoACD dataset, constructed from Korean adolescents’ utterances, contains shorter and more direct statements. These are often marked by explicit modal expressions and strong emotional intensity, reflecting the cognitive style of adolescents, which tends to favor absolutist reasoning and performance-related anxiety. Utterances are typically centered around school life, peer evaluation, and identity development, with minimal narrative elaboration.

In contrast, the Therapist QA dataset consists of English-language clinical dialogues involving adults. These utterances are generally longer, more reflective, and embedded within broader therapeutic narratives. Statements labeled as “Should Statements” frequently appear alongside descriptions of past experiences, diagnoses, and introspective reasoning. This structural complexity introduces interpretive ambiguity that may challenge classification.

These corpus-level differences—in utterance length, discourse structure, and psychological framing—help explain the type-level performance variation observed across datasets. Representative examples from each dataset are provided in Table 15.

LLM	Cognitive Distortion Type	Relevant Sentence (Korean)	Saliency Score
Utterance: <i>I did well in the club interview, but I feel crushed by the thought that I must live up to everyone's expectations. Like the president said, if I don't get selected, does that mean I'm not qualified to be a leader? Do I have to be perfect at everything? I'm definitely lacking. I must do better.</i> (동아리 면접을 봤다. 발표는 잘했지만, 모두의 기대에 부응해야 한다는 생각에 짓눌린다. 회장 언니 말처럼 뽑히지 못하면 나는 리더가 될 자격이 없는 사람일까? 모든 걸 완벽하게 해내야만 하는 건가? 내가 부족한 건 분명하다. 나는 더 잘해야만 해.)			
GPT-4o	Should Statements	<i>I feel crushed by the thought that I must live up to everyone's expectations. ... I must do better.</i> (모두의 기대에 부응해야 한다는 생각에 짓눌린다. ... 나는 더 잘해야만 해.)	0.444
	Labeling	<i>If I don't get selected, does that mean I'm not qualified to be a leader? I'm definitely lacking.</i> (뽑히지 못하면 나는 리더가 될 자격이 없는 사람일까? 내가 부족한 건 분명하다.)	0.333
	Jumping to Conclusions	<i>If I don't get selected, does that mean I'm not qualified to be a leader?</i> (뽑히지 못하면 나는 리더가 될 자격이 없는 사람일까?)	0.222
Gemini	Jumping to Conclusions	<i>I'm definitely lacking.</i> (내가 부족한 건 분명하다.)	0.138
	Emotional Reasoning	<i>I feel crushed by the thought that I must meet everyone's expectations.</i> (모두의 기대에 부응해야 한다는 생각에 짓눌린다.)	0.138
	Should Statements	<i>I must meet everyone's expectations. ... I must do better.</i> (모두의 기대에 부응해야 한다... 나는 더 잘해야만 해.)	0.276
	All-or-Nothing Thinking	<i>If I don't get selected, does that mean I'm not qualified to be a leader?</i> (뽑히지 못하면 나는 리더가 될 자격이 없는 사람일까?)	0.241
	Labeling	<i>If I don't get selected, does that mean I'm not qualified? I'm definitely lacking.</i> (뽑히지 못하면 나는 리더가 될 자격이 없는 사람일까? 내가 부족한 건 분명하다.)	0.207
Claude	Discounting the Positive	<i>I did well in the presentation, but...</i> (발표는 잘했지만)	0.187
	Labeling	<i>I'm definitely lacking.</i> (내가 부족한 건 분명하다.)	0.172
	Overgeneralization	<i>If I don't get selected, does that mean I'm not qualified to be a leader?</i> (뽑히지 못하면 나는 리더가 될 자격이 없는 사람일까?)	0.141
	Should Statements	<i>I feel crushed by the thought that I must live up to everyone's expectations. ... I must do better.</i> (모두의 기대에 부응해야 한다는 생각에 짓눌린다. ... 나는 더 잘해야만 해.)	0.266
	All-or-Nothing Thinking	<i>If I don't get selected, does that mean I'm not qualified to be a leader?</i> (뽑히지 못하면 나는 리더가 될 자격이 없는 사람일까?)	0.234

Table 14: Example of LLM-Inferred Instances from a Single Utterance

Dataset	Utterance
KoACD	<p><i>I did well in the club interview, but I feel crushed by the thought that I must live up to everyone's expectations. Like the president said, if I don't get selected, does that mean I'm not qualified to be a leader? Do I have to be perfect at everything? I'm definitely lacking. I must do better.</i></p> <p>(동아리 면접을 봤다. 발표는 잘했지만, 모두의 기대에 부응해야 한다는 생각에 짓 눌린다. 회장 언니 말처럼 뽑히지 못하면 나는 리더가 될 자격이 없는 사람일까? 모든 걸 완벽하게 해내야만 하는 건가? 내가 부족한 건 분명하다. 나는 더 잘해야만 해.)</p>
Therapist QA	<p><i>I have been suffering from bulimia for four months now. I realize the health risks and I know I have a problem. I have been trying to stop for a month now with no success. Before this problem I was healthy and now I fear that all my hard work I have completed over the years to be a healthy person are going down the drain. To be honest I am not sure what started my ED, but my main focus is to overcome it. I know that I have some self esteem issues and I will continue to work on that, but do you have any advice or tricks to stop these behaviors that have seemed to become habitual and uncontrollable. I know that getting professional help is probably the best way to go, but that is not me. I have always dealt with my problems in the past and I would like to give this a shot. So if you have any suggestions or tips to help me slowly stop these bulimic behaviors I would appreciate it so much.</i></p>

Table 15: Examples of “Should Statements” from KoACD and Therapist QA datasets

F Prompt Templates

We provide the prompt templates used in the LLM-based ELB extraction and cognitive distortion inference processes. The template for ELB extraction is shown in Table 16, and the templates for cognitive distortion instance inference are presented in Table 17 and Table 18.

Emotion–Logic–Behavior Extraction

The user said the following sentence:

"{sentence}"

Please analyze the sentence according to the following three aspects:

1. **Emotion:** Identify emotional elements or affective states expressed or implied in the sentence, such as anger, sadness, anxiety, frustration, or joy.
2. **Logic:** Identify the reasoning or thought process in the sentence. Look for any assumptions, conclusions, generalizations, or causal relationships. Evaluate whether the reasoning is logical or contains any fallacies.
3. **Behavior:** Identify any behaviors or behavioral intentions mentioned in the sentence. Determine whether the person did something, intends to act, or is avoiding action.

For each aspect, provide a brief explanation in the form of a single sentence summarizing the key point.

Please respond in the following JSON format:

```
{  
  "emotion": "One-sentence summary of the emotional aspect",  
  "logic": "One-sentence summary of the logical aspect",  
  "behavior": "One-sentence summary of the behavioral aspect"  
}
```

Note: Each entry should be no more than one sentence. If an aspect is not present, respond with "Not applicable" or "No relevant content found in the sentence."

Table 16: Prompt for Emotion–Logic–Behavior Extraction

Inference of Cognitive Distortion Instances (KoACD)

The user said the following sentence:

“{sentence}”

{emotion_info}

{logic_info}

{behavior_info}

Refer to the definitions of the following 10 cognitive distortion types and output all distortion types and salience scores that can be identified in the sentence above. You must select only from the following 10 types: All-or-Nothing Thinking, Overgeneralization, Mental Filter, Discounting the Positive, Jumping to Conclusions, Magnification and Minimization, Emotional Reasoning, Should Statements, Labeling, Personalization. Do not include any other types.

Identify all cognitive distortions present in the sentence using the predefined types described in Appendix C, as shown in Table 12.

Return all detected cognitive distortions in the following format:

```
{
  {"type": "Cognitive distortion type (must be chosen from the 10 types above)", "salience score": float,
  "relevant_text": "Relevant excerpt from the sentence"},
  {"type": "Cognitive distortion type (must be chosen from the 10 types above)", "salience score": float,
  "relevant_text": "Relevant excerpt from the sentence"},
  ...
}
```

Instructions:

- Cognitive distortion types must be written in English only, exactly as listed above (no Korean or parenthetical explanations).

Example: "All-or-Nothing Thinking" (OK), "All-or-Nothing Thinking" (Not OK)

- Choose only from the following 10 types: All-or-Nothing Thinking, Overgeneralization, Mental Filter, Discounting the Positive, Jumping to Conclusions, Magnification and Minimization, Emotional Reasoning, Should Statements, Labeling, Personalization.

- For each distortion, extract a short and relevant excerpt from the sentence that clearly reflects the distortion.

Strict output format requirements:

1. Return only a JSON array ([]), with no explanations or additional JSON objects.
 2. The array must include all detected cognitive distortion objects.
 3. Each object must contain exactly three fields: "type", "salience score", and "relevant_text".
-

Table 17: Prompt for Inference of Cognitive Distortion Instances (KoACD)

Inference of Cognitive Distortion Instances (Therapist QA)

The user said the following sentence:

“{sentence}”

{emotion_info}

{logic_info}

{behavior_info}

Refer to the definitions of the following 10 cognitive distortion types and output all distortion types and salience scores that can be identified in the sentence above. You must select only from the following 10 types: All-or-nothing thinking, Overgeneralization, Mental filter, Emotional reasoning, Labeling, Magnification, Should statements, Fortune-telling, Mind Reading, and Personalization. Do not include any other types.

Identify all cognitive distortions present in the sentence using the predefined types described in Appendix C, as shown in Table 13.

Return all detected cognitive distortions in the following format:

```
{
  {"type": "Cognitive distortion type (must be chosen from the 10 types above)", "salience score": float,
  "relevant_text": "Relevant excerpt from the sentence"},
  {"type": "Cognitive distortion type (must be chosen from the 10 types above)", "salience score": float,
  "relevant_text": "Relevant excerpt from the sentence"},
  ...
}
```

Instructions:

- Cognitive distortion types must be written in English name only.

Example: "All-or-nothing thinking" (Correct), "All-or-nothing thinking (black and white thinking)" (Incorrect)

- Cognitive distortion types must be selected ONLY from the 10 types provided: All-or-nothing thinking, Overgeneralization, Mental filter, Emotional reasoning, Labeling, Magnification, Should statements, Fortune-telling, Mind Reading, and Personalization. Do not use any other types.

- For each distortion, extract a short and relevant excerpt from the sentence that clearly reflects the distortion.

Strict output format requirements:

1. Return only a JSON array ([]), with no explanations or additional JSON objects.
 2. The array must include all detected cognitive distortion objects.
 3. Each object must contain exactly three fields: "type", "salience score", and "relevant_text".
-

Table 18: Prompt for Inference of Cognitive Distortion Instances (Therapist QA)

G Evaluation Form

Expert Evaluation Form for Cognitive Distortion

Korean

안녕하세요, 선생님.

본 연구에 참여해주셔서 진심으로 감사드립니다.

이 작업은 청소년의 실제 발화를 기반으로, 해당 발화에 내포된 인지 왜곡 유형을 평가해주는 작업입니다.

데이터는 파일에 준비되어 있습니다.

각 발화에 대해 아래의 기준에 따라 가장 적절한 인지 왜곡 유형을 선택한 후, 해당 유형에 대한 신뢰도 점수를 1-3점 척도로 평가해 주시면 됩니다.

1점: 선택된 유형이 문맥에 부합하지 않거나 명확하지 않음

2점: 해당 유형이 어느 정도 나타나지만, 다른 유형과 혼동의 여지가 있음

3점: 해당 인지 왜곡 유형이 명확하고 주된 특징으로 드러남

또한, 해당 발화가 다른 인지 왜곡 유형으로도 해석될 가능성이 있다고 판단되시는 경우, 구체적으로 어떤 유형으로도 해석될 수 있는지와 그에 대한 간단한 이유를 함께 적어주시면 감사하겠습니다.

작업 중 불편함이나 감정적 부담을 느끼시는 경우, 언제든지 작업을 중단하거나 일정 조정을 요청하실 수 있습니다.

문의사항이 있으실 경우 언제든지 연락 주시기 바랍니다.

감사합니다.

English

Hello.

Thank you very much for participating in this study.

This task is to evaluate the types of cognitive distortions contained in actual utterances of adolescents.

The data is prepared in the file.

For each utterance, select the most appropriate type of cognitive distortion according to the criteria below, and then evaluate the reliability score on a scale of 1 to 3.

1 point: The selected type does not fit the context or is unclear.

2 points: The type is somewhat present but may be confused with other types.

3 points: The cognitive distortion type is clear and stands out as the main feature.

Additionally, if you believe that the utterance could be interpreted as another type of cognitive distortion, Please specify which other types it could be interpreted as and provide a brief explanation of why.

If you experience any discomfort or emotional distress during the process, you may stop at any time or request a schedule adjustment.

If you have any questions, please feel free to contact us at any time.

Thank you.

Table 19: Expert Evaluation Form for Cognitive Distortion