

CareCorpus+ for Environment: Extended Classification and Augmentation of Caregiver Strategies to Create Enabling Environments for Young Children with Disabilities

Anonymous ACL submission

Abstract

Caregivers play a central role in supporting the participation of children with disabilities in everyday activities. Environmental strategies, such as rearranging routines or spaces, are central to this support, yet they remain scarce and difficult to model systematically due to their free-text nature. To address this gap, we present 1,848 environment-focused caregiver strategies from pediatric rehabilitation contexts, manually annotated into five clinically grounded sub-categories capturing activity demands, social contexts, physical environments, policies, and resources. Using this dataset, we benchmark a multi-class classification task and evaluate both manual and synthetic data augmentation. BERT-based models achieve improvements of up to 32% in macro-F₁, supporting systematic evaluation of NLP methods for modeling caregiver strategies in low-resource rehabilitation settings.

1 Introduction

Participation, defined as children’s attendance and involvement in everyday activities across different settings, is a central outcome in pediatric rehabilitation (Imms et al., 2017, 2016; World Health Organization, 2007). Pediatric rehabilitation aims to support children’s engagement in daily life, making participation a key indicator of whether interventions translate into meaningful, real-world impact. However, achieving meaningful participation requires more than children’s individual abilities, as it is also shaped by the environmental factors that surround them.

Prior work consistently shows that environmental factors, including the physical, social, and organizational contexts, are key drivers of participation in valued activities (Albrecht and Khetani, 2017; Di Marino et al., 2018; Khetani et al., 2018b, 2020, 2018a). Caregiver-report studies further show that caregivers most frequently describe environmental

strategies, such as adjusting routines or modifying physical spaces, as contributing to improvements in their child’s participation (Kaelin et al., 2021; Jarvis et al., 2019b). These environmental strategies are most often described in free-text narratives, which makes it challenging to systematically organize and leverage them to help caregivers identify appropriate support options.

As a result, recent NLP work has begun to explore automated methods for detecting and classifying caregiver-generated strategies (Valizadeh et al., 2024; Farzana et al., 2024; Kaelin et al., 2021) often drawing upon data from a suite of web-based tools known as the Participation and Environment Measure (PEM) (Khetani et al., 2015; Coster et al., 2011; Shahin et al., 2024). However, this existing work (Valizadeh et al., 2024; Farzana et al., 2024; Kaelin et al., 2021) groups environmental strategies under a single category, despite caregivers describing a wide range of distinct environmental supports. This coarse representation, limited partly due to the low-resource nature of caregiver-generated data, constrains understanding of how participation is enabled in different settings. We address this limitation by making the following contributions:

- We introduce **fine-grained annotations of environmental caregiver strategies**, applying clinically grounded subcategories to capture how caregivers modify their child’s environment to support participation.
- We formulate and evaluate a fine-grained classification task using these strategies, **providing empirical baselines** that demonstrate the feasibility of automatic classification from caregiver-authored text.
- We show that both manual and synthetic augmentation of caregiver strategies improve fine-grained classification performance, **highlighting augmentation as a practical approach**

081	for low-resource caregiver strategy modeling.	low-resource conditions. Within a pediatric reha-	131
082		ilitation context, data augmentation has proven	132
083	The remainder of the paper describes the dataset	feasible at a coarse-grained level (Farzana et al.,	133
084	and annotation scheme, followed by task formula-	2024). Our study builds on this line of work, but dif-	134
085	tion, augmentation methods, experimental results,	fers (1) by focusing on caregiver strategies within	135
	and discussion.	the environmental class given its central role as	136
086	2 Related Work	a determinant of participation outcomes and (2)	137
087	Prior work has classified caregiver strategies into	by evaluating the viability of multiple augmenta-	138
088	four classes in rehabilitation (<i>Environment/Context,</i>	tion approaches for responsibly approaching the	139
089	<i>Sense of Self, Preferences,</i> and <i>Activity Compe-</i>	development of downstream AI-enabled pediatric	140
090	<i>tence</i>) and a “no-strategy” class for irrelevant strate-	rehabilitation applications.	141
091	gies (Farzana et al., 2024; Valizadeh et al., 2024).		
092	However, it has not further separated the Environ-	3 Data Annotation	142
093	ment/Context class into its sub-components. Ev-		
094	idence suggests that caregivers most commonly	3.1 Data Source	143
095	target the environment when supporting children’s		
096	participation across settings (Jarvis et al., 2019b;	Our data are derived from the CareCorpus+ dataset	144
097	Kaelin et al., 2021; Khetani and Lucero, 2024), mo-	(Farzana et al., 2024), a large corpus of 3,062	145
098	tivating the need to examine environmental strate-	caregiver-generated strategies collected through the	146
099	gies beyond a single coarse category to support	Young Children’s Participation and Environment	147
100	more effective and individualized participation-	Measure (YC-PEM) (Khetani et al., 2015). Each	148
101	focused care planning.	strategy describes an action that caregivers take	149
102	Few existing datasets are suitable for classifying	(or could take) to help their child participate in	150
103	environmentally-focused caregiver strategies in	daily activities across home, daycare/preschool,	151
104	pediatric rehabilitation settings. Prior clinical or	and community settings. Strategies were manu-	152
105	behavioral datasets do not contain indicators of en-	ally categorized into five clinically grounded con-	153
106	vironmental supports, and have instead linked free-	structs: Environment/Context, Sense of Self, Ac-	154
107	text clinical documents to ICF codes using auto-	tivity Competence, Preferences, and Non-Strategy,	155
108	mated approaches (Kukafka et al., 2006; Newman-	which align with established pediatric rehabilita-	156
109	Griffis et al., 2021) and collected video-recorded	tion frameworks for child participation (Imms et al.,	157
110	sessions of children with ASD and typically devel-	2017; Kaelin et al., 2023; Valizadeh et al., 2024).	158
111	oping children to study engagement (Chori-	In this study, we focus exclusively on the Envi-	159
112	anopoulou et al., 2017). This gap is particularly	ronment/Context subset, the largest and most di-	160
113	consequential for participation-focused pediatric	verse class within CareCorpus+ ($n=1,848$). This	161
114	rehabilitation, where caregivers use environmental	subset captures environmental supports and contex-	162
115	strategies to support children’s involvement in ev-	tual factors that caregivers describe when facilitat-	163
116	eryday life contexts (Jarvis et al., 2019b; Kaelin	ing their child’s participation in valued activities	164
117	et al., 2021; Khetani and Lucero, 2024). The clos-	across different service contexts (e.g., home rou-	165
118	est related dataset is CareCorpus+ (Farzana et al.,	tines, community outings, daycare settings) and	166
119	2024), which contains 3,062 caregiver strategies	reflects how these conditions can be adjusted or	167
120	organized into categories aligned with established	leveraged to support engagement.	168
121	participation constructs (Imms et al., 2017), one		
122	of which is Environment/Context. It includes care-	3.2 Annotation Process	169
123	giver strategy data collected during early interven-		
124	tion trials of the PEM for children across broad	We sought to further categorize caregiver strategies	170
125	age ranges (0-5) and diverse pediatric rehabilita-	within the environmental class into finer-grained	171
126	tion contexts (Kaelin et al., 2023; Valizadeh et al.,	subcategories that capture distinct approaches to	172
127	2024; Farzana et al., 2024).	targeting the child’s environment to improve their	173
128	Pediatric healthcare settings often suffer from	participation in activities. This refinement enables	174
129	small and imbalanced clinical data, making data	the identification of specific ways that caregivers	175
130	augmentation a valuable strategy for addressing	modify the young child’s environment to support	176
		their participation. It also supports the develop-	177
		ment of automated strategy classification models	178
		for integration into family-centered care-planning	179

180 applications.

181 Annotations were conducted by two trained un-
182 dergraduate research assistants majoring in com-
183 puter science and neuroscience, as supervised by a
184 research-engaged pediatric occupational therapist
185 and a computer science PhD student. The anno-
186 tation guideline, inspired by Farzana et al. (2024)
187 and developed under the guidance of a rehabilita-
188 tion scientist, outlined five Environment/Context
189 subtypes, each corresponding to a label (1-5):

- 190 • **1: Demands of the Activity:** Represents the
191 physical (e.g., walking), cognitive (e.g., at-
192 tention), or social (e.g., taking turns) effort
193 required for a child to engage in an activity.
- 194 • **2: Relationships / Attitudes:** Captures strate-
195 gies emphasizing the support, participation,
196 or attitudes of other people (e.g., "Mom sings
197 while child brushes teeth.").
- 198 • **3: Physical / Sensory:** Relates to the physical
199 setup, layout, or sensory qualities of the envi-
200 ronment (e.g., "Rearranged furniture to create
201 more play space.").
- 202 • **4: Weather / Safety / Policies:** Encompasses
203 rules, environmental conditions, and safety
204 considerations that shape participation (e.g.,
205 "Avoid crossing busy streets.").
- 206 • **5: Resources:** Refers to tangible or intangible
207 supports that make participation possible (e.g.,
208 "Therapist visits on Mondays.").

209 If a strategy did not fit into any of the above cat-
210 egories, annotators were told to assign the additional
211 label, **0: Misclassified**. The annotators indepen-
212 dently labeled a total of 1,848 strategies over seven
213 weeks, averaging around 260 strategies per week.
214 After each annotation round, they met with the su-
215 pervising researchers, including a domain expert, to
216 discuss and resolve any disagreements. This itera-
217 tive adjudication process clarified ambiguous cases
218 and aligned interpretation between annotators. Fig-
219 ure 1 summarizes the final label distribution.

220 3.3 Inter-Annotator Agreement

221 To ensure consistency, we computed agreement
222 for all labels, as shown in Table 1. Agreement
223 values ranged from 64.8% to 83.9%, with an overall
224 agreement of 72.8%. We also computed Cohen’s
225 κ (Cohen, 1960) to complement the overall score,
226 which yielded a substantial agreement of $\kappa = 0.72$.

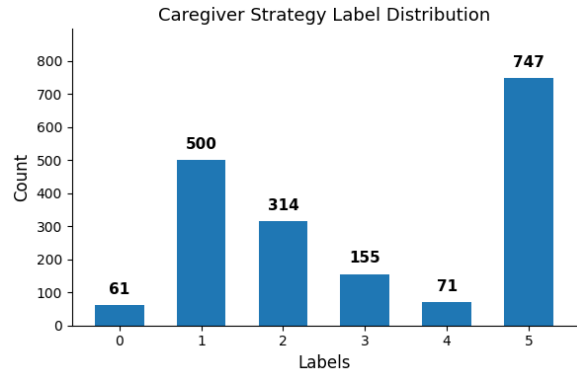


Figure 1: Distribution of annotated caregiver strat-
egy labels. The numbers above each bar indicate
the count of strategies per category. Labels cor-
respond to: 0–Misclassified, 1–Demands of the Activ-
ity, 2–Relationships/Attitudes, 3–Physical/Sensory, 4–
Weather/Safety/Policy, and 5–Resources.

Label (#)	Agreement (%)
Demands of the Activity (1)	79.5
Relationships/Attitudes (2)	83.9
Physical/Sensory (3)	70.4
Weather/Safety/Policy (4)	65.7
Resources (5)	64.8
Overall	72.8

Table 1: Inter-annotator agreement by label. Agree-
ment (%) reflects the proportion of annotator pairs who
assigned the same label. Overall Agreement was com-
puted as the class size-weighted average across all rows.

227 3.4 Dataset Summary

228 Among the five annotated subtypes, *Resources*
229 ($n=747$ strategies) and *Demands of the Activity*
230 ($n=500$) were the most prevalent, indicating that
231 caregivers actively leverage tangible supports, ser-
232 vices, or activity-related demands to facilitate par-
233 ticipation. *Weather/Safety/Policy* strategies were
234 relatively rare ($n=71$), suggesting that environ-
235 mental constraints or enforced rules were less fre-
236 quently emphasized. Linguistic patterns also varied
237 across labels. For example, in *Resources* (Figure
238 2), frequent bigrams included “age appropriate,”
239 “make sure,” and “having time,” showing how care-
240 givers secure or plan supports such as materials,
241 time, or information. Meanwhile, *Demands of*
242 *the Activity* strategies often contained verbs like
243 “show,” “practice,” or “read,” describing modeling
244 and skill-building approaches. Overall, the distri-
245 bution illustrates that caregivers tend to adopt con-
246 crete, resource-oriented strategies in daily routines.

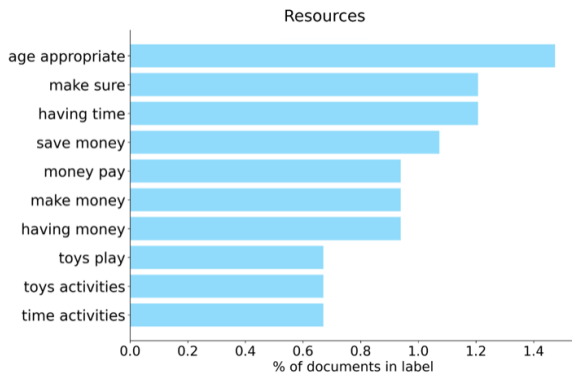


Figure 2: Top 10 most frequent bigrams in *Resources* strategies.

Additional bigrams are provided in Appendix A.

4 Data Augmentation

Research in pediatric rehabilitation often relies on caregiver-generated data that must be manually elicited and annotated, making large-scale data collection resource-intensive. This is particularly challenging for fine-grained analysis, where labeled data for individual strategy subcategories are limited. Data augmentation therefore represents a potentially valuable approach for scaling the training data.

4.1 Manual Data Augmentation

To investigate the feasibility of augmenting environmentally focused caregiver strategies, we developed a manual augmentation protocol that draws on commonly referenced text augmentation categories described by Wei and Zou (2019), adapted for a human-guided approach. The protocol includes four rewriting techniques: Selective Deletion (D), Selective Replacement (R), Selective Addition (A) and Word Shuffling (S). Descriptions and representative examples of each technique are provided in Table 2, with the full version available in Appendix A. The manual augmentation process was initially piloted by two individuals: one with a neuroscience background and one with a computer science background. The individuals were instructed to focus their augmentations on two underrepresented classes (*Physical/Sensory* and *Weather/Safety/Policy*) identified in Figure 1. The pilot phase was included to guide the protocol and to identify how different forms of manual augmentation could be implemented in the main study.

After the pilot phase, a pediatric rehabilitation

clinical expert manually augmented 1,786 strategies following the same protocol. Among the four techniques, Selective Replacement (R) was used most frequently ($n=1,732$; 41.9%), whereas Selective Addition (A) was used least frequently ($n=118$; 2.9%). These patterns were consistent in home, daycare/preschool, and community settings.

4.2 Synthetic Data Augmentation

To further explore whether large language models can effectively generate synthetic caregiver strategies, we employed GPT-5.1 (OpenAI, 2025) to rewrite training examples that mirror the human augmentation process. We implemented synthetic strategy generation through a two-step process, involving an initial strategy generation step followed by an evaluation and selection step to refine the set of strategies included in the augmented output.

4.2.1 Strategy Generation

Each caregiver strategy was rewritten using four different rewriting styles drawn from the human augmentation guideline: Selective Deletion (D), Selective Replacement (R), Selective Addition (A), and Word Shuffling (S). Each method was prompted separately to prevent instruction interference and ensure higher-quality outputs, following the idea of task decomposition by Khot et al. (2022). The prompt formats for each rewriting style are shown in Table 2. Each prompt begins with a common instruction asking the model to rewrite the original strategy according to a specified method, followed by the method’s definition and a real example from the dataset to guide the model. This process was iterated for all training-set strategies (excluding label 5 which is *Misclassified/not an EC strategy*), yielding 5,716 synthetic strategies. Table 3 provides manually and synthetically generated examples for each caregiver strategy category.

4.2.2 Evaluation and Selection

Each synthetically rewritten version across D, R, A, and S was evaluated using three combined criteria: semantic similarity (sim), fluency (perplexity, ppl), and GPT-based self-evaluation score (gpt, 1–5 scale). This evaluation setup was designed to mirror the human selection process used during manual augmentation, where annotators selected the most appropriate rewritten strategy based on meaning preservation and naturalness. Semantic similarity measures how well the core meaning of the original caregiver strategy is preserved, while

Shared Instruction: Rewrite the following caregiver strategy sentence using the specified method while preserving the original meaning.

Approach	Augmentation Guideline
D	Remove unnecessary words or verbs to isolate the main action or object. Example: "Rocking him while singing a song and then putting him down awake" → "Rocking while singing."
R	Replace a key adjective or verb with a natural synonym that preserves meaning. Example: "Shut off TV" → "Turn off TV."
A	Add one realistic, relevant action or detail that complements the strategy. Example: "Wheelchair accessibility" → "Ensure wheelchair accessibility."
S	Reorder clauses or phrases while keeping the same meaning. Example: "When home alone allow more freedom" → "Allow more freedom when home alone."

Table 2: Summary of the four augmentation approaches used for both manual and synthetic data generation. Approaches are **D**=Selective Deletion, **R**=Selective Replacement, **A**=Selective Addition, and **S**=Word Shuffling. The definitions were adapted from the human annotation guideline and used directly as prompt instructions for GPT-based synthetic strategy generation. See the full human annotation guideline in Appendix A.

fluency, captured through perplexity, reflects the readability and naturalness of the generated output. The GPT-based self-evaluation allows the model to make its own judgment by scoring each rewritten version on a 1–5 scale using the prompt:

Rate each version on a scale from 1 to 5 based on how well it preserves the original meaning while sounding natural and fluent.

- 1 = Very poor (meaning changed or sentence unnatural)
- 2 = Weak (partially preserves meaning but awkward phrasing)
- 3 = Fair (mostly correct meaning but minor awkwardness)
- 4 = Good (meaning well-preserved, natural phrasing)
- 5 = Excellent (fully preserves meaning, very natural and fluent)

Considering semantic, linguistic, and model-based self-evaluation together, our framework holistically assesses the overall quality of each rewritten strategy. We normalized all metrics to the [0,1] range and computed a weighted composite score for each version v as:

$$\text{score}_v = \alpha \cdot \text{sim}_v^{\text{norm}} + \beta \cdot (1 - \text{ppl}_v^{\text{norm}}) + \gamma \cdot \text{gpt}_v^{\text{norm}} \quad (1)$$

where $\alpha = 0.5$, $\beta = 0.3$, and $\gamma = 0.2$ represent the respective weights for semantic similarity, fluency, and GPT-based self-evaluation. We chose these weights to prioritize semantic preservation, followed by fluency, while treating GPT-based self-evaluation as a lower-weight auxiliary signal due to potential bias. The version with the highest composite score across D , R , A , and S was selected as the final synthetic rewrite, thereby narrowing the synthetic examples to 1,429.

Across all synthetic strategies generated, **R** was the most frequently chosen method (909 strategies), followed by **D** (396), **A** (119), and **S** (5). Based on the combined evaluation criteria, replacement-based rewrites appeared to have the best balance of meaning preservation and fluency within our framework. This trend is broadly similar to the distribution observed in the human-augmented dataset, where **R** and **D** were most commonly selected.

5 Strategy Classification

After the data annotation and augmentation stages, we assessed the dataset’s utility through the established benchmark task of automated caregiver strategy classification (Kaelin et al., 2023; Valizadeh et al., 2024; Farzana et al., 2024). This helped clarify the capacity to automatically model fine-grained Environment/Context strategy distinctions, thereby establishing the dataset’s applicability for downstream rehabilitation contexts.

Category	Strategy
Demands of the Activity	Original: "Show how to do the activity and allow him to copy." → Manual: "Rocking while singing." (<i>Deletion</i>) → Synthetic: "Demonstrate how to do the activity and permit him to copy." (<i>Replacement</i>)
Relationships/Attitudes	Original: "Sing song together." → Manual: "Chant songs together." (<i>Replacement</i>) → Synthetic: "Sing songs together when needed." (<i>Addition</i>)
Physical/Sensory	Original: "Make easy for him to access his toys/clothes." → Manual: "Easy access to toys/clothes." (<i>Deletion</i>) → Synthetic: "Create easy for him to access his toys/clothes." (<i>Replacement</i>)
Weather/Safety/Policy	Original: "Control the amount of TV or tablet time to make time for these other activities." → Manual: "Control the TV or tablet time for other activities." (<i>Deletion</i>) → Synthetic: "Control the amount of TV or tablet time to create time for these other activities." (<i>Replacement</i>)
Resources	Original: "Exercise programs - for young children to help them keep their senses in balance." → Manual: "Help young children keep their senses in balance through exercise programs." (<i>Substitution</i>) → Synthetic: "Exercise programs - for young children to help them keep their senses in balance to support development." (<i>Addition</i>)

Table 3: Examples of manual and synthetic augmentations for each caregiver strategy category. Each rewrite is followed by its corresponding augmentation type (e.g., Deletion, Replacement, Addition, or Substitution).

5.1 Experimental Setup

We framed our task as a multi-class text classification problem, where the goal was to predict the correct strategy label for each caregiver strategy. The dataset was divided into training (80%) and test (20%) sets using a stratified split to preserve the original class distribution. Model performance was evaluated on the fixed test set using macro-averaged accuracy, precision, recall, and F₁ scores. We employed simple baselines and transformer-based models to benchmark performance:

- **Most Frequent Class:** A simple baseline that predicts the most frequent label observed in the training data. This serves as a performance floor that reflects class imbalance.
- **Random:** A stochastic baseline that assigns labels uniformly at random across all classes, representing performance without any learned information.
- **Naive Bayes:** A probabilistic model based on word frequency distributions with a multinomial likelihood, assuming conditional independence among features given the predicted label.
- **Logistic Regression:** A linear classifier

trained on TF-IDF feature representations of the input text, modeling the probability of each class through a softmax layer.

- **BERT:** A transformer-based model (Devlin et al., 2019) fine-tuned on the training data to learn contextual embeddings of the input text. We used *bert-base-uncased* from HuggingFace¹ with a linear classification head, trained for three epochs with a batch size of 16, a learning rate of $2e-5$.
- **BioClinicalBERT:** A domain-adapted variant of BERT pre-trained on biomedical and clinical corpora (Alsentzer et al., 2019). The model was fine-tuned on the training data under the same setup as BERT to examine whether domain-specific pre-training improves contextual understanding of the input caregiver strategy text.
- **LLaMA:** A large instruction-tuned language model (Touvron et al., 2023) evaluated without task-specific fine-tuning. We tested both zero-shot and few-shot prompting conditions, using two representative examples per class selected based on their proximity to the class centroid in TF-IDF space. The model was not

¹<https://huggingface.co/docs/transformers>

436 fine-tuned on augmented datasets to keep the
437 evaluation consistent.

438 We compared strategy classification perfor-
439 mance across baselines and transformer-based mod-
440 els under three data configurations: the original
441 annotated dataset, the original dataset combined
442 with manually augmented strategies, and the origi-
443 nal dataset combined with synthetically generated
444 strategies. This setup allowed us to observe how
445 different types of augmentation affect model per-
446 formance. All fine-tuning experiments were con-
447 ducted on a single NVIDIA A100 GPU.

448 5.2 Results

449 Model performance across all data configurations
450 is summarized in Table 4. Results for BERT and
451 BioClinicalBERT are averaged across three runs,
452 and macro-averaged metrics are reported to account
453 for class imbalance across strategy types.

454 Overall, BERT consistently performed the best
455 across all datasets, aligning with prior findings
456 from Farzana et al. (2024) in coarser-grained tasks.
457 Its performance further improved after augmen-
458 tation, reaching the best overall performance of
459 $F_1=0.863$ on the *Original + Manual Augmenta-*
460 *tion* dataset. BioClinicalBERT followed closely,
461 showing consistently strong results across datasets
462 and achieving a matching $F_1=0.739$ on the *Orig-*
463 *inal + Synthetic Augmentation* dataset. Logis-
464 tic Regression and Naive Bayes achieved moder-
465 ate performance, ranging from 0.422 to 0.668 F_1 ,
466 while the Most Frequent and Random baselines
467 remained lowest as expected. Prompting LLaMA
468 performed comparably to Logistic Regression and
469 Naive Bayes, with a zero-shot $F_1=0.489$ and a
470 slightly improved few-shot $F_1=0.591$.

471 Augmenting the data generally improved clas-
472 sification performance, especially for transformer-
473 based models. After manual augmentation (*Orig-*
474 *inal + Manual Aug*), the top-performing model,
475 BERT, improved from an F_1 score of 0.656 to
476 0.863, an increase of approximately 32%. Bio-
477 ClinicalBERT exhibited a similar trend (0.632
478 \rightarrow 0.830, an increase of 31%), showing that
479 domain-specific pretraining also benefits from ad-
480 ditional data. Naive Bayes and Logistic Regression
481 achieved modest gains (roughly 10–15%).

482 After synthetic augmentation (*Original + Syn-*
483 *thetic Aug*), Naive Bayes showed a minor decrease
484 (0.424 \rightarrow 0.422 in F_1) compared to the *Original*
485 dataset, but all other models improved in a pat-

486 tern consistent with manual augmentation. Al-
487 though the overall gains were smaller in *Orig-*
488 *inal + Synthetic Aug*, BERT and BioClinicalBERT
489 both reached an F_1 score of 0.739, which is ap-
490 proximately 13%-17% higher than their original
491 performance. These results suggest that while
492 synthetically-generated rewrites do not yet match
493 the quality of human-generated ones, they mean-
494 ingfully enhance model robustness.

495 6 Error Analysis

496 We conducted a qualitative error analysis on BERT
497 fine-tuned on the *Original + Synthetic Augmenta-*
498 *tion* dataset, the best-performing GPT-augmented
499 model, to identify where synthetic data contributed
500 to misclassification and how future augmentation
501 could be refined. The analysis was performed on
502 the fixed test set (358 samples), of which 70 (ap-
503 proximately 19.5%) were misclassified.

504 Most confusions occurred between *Demands*
505 *of the activity* \leftrightarrow *Resources* (25 cases) and *De-*
506 *mands of the activity* \leftrightarrow *Relationships/attitudes* (10
507 cases). Many of these cases were also challenging
508 for human annotators to separate during the initial
509 data annotation phase. For instance, “Watching a
510 YouTube video or reading a book on the topic” was
511 labeled *Demands of the activity* but predicted as *Re-*
512 *sources*, likely because of its references to tangible
513 objects (e.g., video and book) that resemble strate-
514 gies in *Resources*. During manual annotation, this
515 example was discussed and labeled *Demands of*
516 *the activity* since the emphasis was placed slightly
517 more on the action than the material itself. Con-
518 versely, “Once something doesn’t work, try a new
519 strategy” was labeled *Resources* but predicted as
520 *Demands of the activity*, as the model prioritized
521 the verb “try” (the action) rather than interpreting
522 the “strategy” as the main resource. These border-
523 line cases reveal the need for explicit tie-breaking
524 criteria in the annotation protocol and for targeted
525 inclusion of ambiguous examples in data augmenta-
526 tion, to help models learn finer distinctions between
527 potentially overlapping strategy categories.

528 6.1 Discussion

529 Our experiments demonstrated that both human-
530 based (manual) and synthetic augmentation im-
531 proved performance across all models, with
532 transformer-based models such as BERT and Bio-
533 ClinicalBERT achieving the highest scores. These
534 results support the feasibility of using data aug-

Data	Model	Accuracy	Precision	Recall	F ₁
Original	Most Frequent	0.419	0.084	0.200	0.118
	Random	0.182	0.184	0.200	0.164
	Naive Bayes	0.616	0.855	0.397	0.424
	Logistic Regression	0.654	0.817	0.465	0.508
	BERT	0.743	0.737	0.626	0.656
	BioClinicalBERT	0.760	0.792	0.609	0.632
	LLaMA (zero-shot)	0.528	0.548	0.511	0.489
	LLaMA (few-shot)	0.643	0.598	0.597	0.591
Original + Manual Aug	Most Frequent	0.418	0.084	0.200	0.118
	Random	0.206	0.205	0.210	0.184
	Naive Bayes	0.700	0.873	0.485	0.529
	Logistic Regression	0.763	0.866	0.607	0.668
	BERT	0.894	0.876	0.853	0.863
	BioClinicalBERT	0.874	0.849	0.815	0.830
Original + Synthetic Aug	Most Frequent	0.419	0.084	0.200	0.118
	Random	0.201	0.216	0.221	0.183
	Naive Bayes	0.617	0.848	0.395	0.422
	Logistic Regression	0.668	0.829	0.478	0.522
	BERT	0.791	0.751	0.728	0.739
	BioClinicalBERT	0.802	0.769	0.717	0.739

Table 4: Performance of models on caregiver strategy classification with different training data variants. “Original” is the original, non-augmented training set; “Original + Manual Aug” adds manually paraphrased examples; and “Original + Synthetic Aug” includes synthetically generated examples. All metrics are macro-averaged to account for class imbalance. LLaMA models were evaluated only on the original dataset, without augmented training, as they serve as large pre-trained baselines rather than fine-tuned models to maintain the same evaluation setup.

535 mentation to expand caregiver strategy datasets, 536 which are inherently low-resource. In our results, 537 BioClinicalBERT performed more comparably to 538 BERT after synthetic augmentation. Qualitative 539 inspection revealed that synthetic strategies often 540 adopted more directive, action-oriented language 541 (e.g., “prepare” or “ensure”). BioClinicalBERT cor- 542 rectly classified examples like “*Ensure wheelchair* 543 *access*” and “*Provide samples of various foods*” 544 which BERT misclassified. This pattern suggests 545 that clinically pretrained models may be better at- 546 tuned to the procedural language caregivers use 547 when describing environmental strategies.

548 From an application perspective, these findings 549 have direct implications for technology-supported 550 caregiver tools, where caregivers often have ac- 551 cess to large banks of strategies but face challenges 552 in browsing or identifying relevant options. Au- 553 tomatic classification and expansion of strategies 554 can facilitate navigation and support individual- 555 ized goal-setting. While manual augmentation 556 yielded higher classification performance, human-

557 generated rewrites could be leveraged to guide and 558 refine future synthetic generation, enabling scal- 559 able and cost-efficient augmentation pipelines for 560 pediatric rehabilitation applications.

561 7 Conclusion

562 We introduce an expanded and carefully annotated 563 caregiver strategy dataset ($n = 1,848$) focused on 564 environmental strategies supporting participation 565 for children with disabilities. We also expand 566 this dataset manually and synthetically, establish- 567 ing guidelines for generating and evaluating aug- 568 mented caregiver strategies. Using this framework, 569 we demonstrate that the augmented data yield per- 570 formance gains of up to 32% in F₁ scores (from 571 0.63–0.66 to 0.74–0.86 for BERT-based models) on 572 fine-grained strategy classification. These results 573 highlight the potential for scaling and applying 574 NLP tools within pediatric rehabilitation, particu- 575 larly to support caregivers in low-resource clinical 576 contexts.

577 Limitations

578 While this work demonstrates the feasibility of
579 augmenting and classifying caregiver strategies,
580 several limitations remain. First, our evaluation
581 was conducted on a relatively small and domain-
582 specific dataset, which may limit generalizability
583 to broader healthcare contexts. Second, although
584 both manual and synthetic augmentation improved
585 performance, the quality of synthetic rewrites was
586 sometimes less robust or natural-sounding com-
587 pared to human-generated ones. Future work could
588 incorporate human-in-the-loop evaluation and ex-
589 plore more prompting techniques to improve the
590 quality of synthetic data.

591 Ethical Considerations

592 This study used de-identified caregiver strategy
593 data from the CareCorpus+ dataset, which was orig-
594 inally collected under institutional ethics approval
595 as described in prior work (Khetani et al., 2015;
596 Kaelin et al., 2022; Khetani et al., 2023; Rizk et al.,
597 2023; Choong et al., 2018; Jarvis et al., 2019a;
598 Khetani et al., 2018a). All participants in the orig-
599 inal studies provided informed consent prior to
600 participation. No identifiable personal or health
601 information was used in this research, and all anal-
602 yses were conducted solely for research purposes
603 to support the development of tools that assist care-
604 givers in browsing and organizing strategies more
605 efficiently. All annotators involved in providing
606 fine-grained Environment/Context strategy labels
607 and manual strategy augmentations completed insti-
608 tutional ethics training. They were compensated for
609 their annotation and augmentation work through re-
610 search assistantships paid at institutionally-defined
611 rates set higher than local minimum wage.

612 References

613 Erin C Albrecht and Mary A Khetani. 2017. Environ-
614 mental impact on young children’s participation in
615 home-based activities. *Developmental Medicine &*
616 *Child Neurology*, 59(4):388–394.

617 Emily Alsentzer, John Murphy, William Boag, Wei-
618 Hung Weng, Di Jindi, Tristan Naumann, and
619 Matthew McDermott. 2019. Publicly available clinical
620 bert embeddings. In *Proceedings of the 2nd clinical*
621 *natural language processing workshop*, pages
622 72–78.

623 Karen Choong, Douglas Fraser, Samah Al-Harbi, Asm
624 Borham, Jill Cameron, Saoirse Cameron, Ji Cheng,

Heather Clark, Tim Doherty, Nora Fayed, and 1 oth- 625
ers. 2018. Functional recovery in critically ill chil- 626
dren, the “weecover” multicenter study. *Pediatric* 627
Critical Care Medicine, 19(2):145–154. 628

Arodami Chorianopoulou, Efthymios Tzinis, Elias Iosif, 629
Asimena Papoulidi, Christina Papailiou, and Alexan- 630
dros Potamianos. 2017. Engagement detection for 631
children with autism spectrum disorder. In *Proceed-* 632
ings of the 2017 IEEE International Conference on 633
Acoustics, Speech and Signal Processing (ICASSP), 634
pages 5055–5059, New Orleans, USA. IEEE. 635

Jacob Cohen. 1960. A coefficient of agreement for 636
nominal scales. *Educational and psychological mea-* 637
surement, 20(1):37–46. 638

Wendy Coster, Gary Bedell, Mary Law, Mary A. 639
Khetani, Rachel Teplicky, Kendra Liljenquist, Kara 640
Gleason, and Ying-Chia Kao. 2011. *Psychometric* 641
evaluation of the participation and environment mea- 642
sure for children and youth. Developmental medicine 643
and child neurology, 53(11):1030–1037. 644

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 645
Kristina Toutanova. 2019. Bert: Pre-training of deep 646
bidirectional transformers for language understand- 647
ing. In *Proceedings of the 2019 conference of the* 648
North American chapter of the association for com- 649
putational linguistics: human language technologies, 650
volume 1 (long and short papers), pages 4171–4186. 651

Erica Di Marino, Stephanie Tremblay, Mary Khetani, 652
and Dana Anaby. 2018. The effect of child, fam- 653
ily and environmental factors on the participation 654
of young children with disabilities. *Disability and* 655
health journal, 11(1):36–42. 656

Shahla Farzana, Ivana Lucero, Vivian Villegas, Vera C 657
Kaelin, Mary Khetani, and Natalie Parde. 2024. 658
Carecorpus+: expanding and augmenting caregiver 659
strategy data to support pediatric rehabilitation. In 660
The 2024 Conference on Empirical Methods in Natu- 661
ral Language Processing (EMNLP), Miami, Florida, 662
USA, November 12-16, 2024., pages 6912–6927. As- 663
sociation for Computational Linguistics. 664

Christine Imms, Brooke Adair, Deb Keen, Anna Ullen- 665
hag, Peter Rosenbaum, and Mats Granlund. 2016. 666
'participation': a systematic review of language, 667
definitions, and constructs used in intervention re- 668
search with children with disabilities. *Developmental* 669
medicine and child neurology. 670

Christine Imms, Mats Granlund, Peter H Wilson, Bert 671
Steenbergen, Peter L Rosenbaum, and Andrew M 672
Gordon. 2017. Participation, both a means and an 673
end: a conceptual analysis of processes and outcomes 674
in childhood disability. *Developmental Medicine &* 675
Child Neurology, 59(1):16–25. 676

Jessica M Jarvis, Karen Choong, and Mary A Khetani. 677
2019a. Associations of participation-focused strate- 678
gies and rehabilitation service use with caregiver 679
stress after pediatric critical illness. *Archives of phys-* 680
ical medicine and rehabilitation, 100(4):703–710. 681

682	Jessica M. Jarvis, Anupama R. Gurga, Hyunji Lim, Jessica Cameron, Jan Willem Gorter, Karen Choong, and Mary A. Khetani. 2019b. Caregiver strategy use to promote children’s home participation after pediatric critical illness . <i>Archives of Physical Medicine and Rehabilitation</i> , 100(11):2144–2150.	738
683		739
684		740
685		741
686		742
687		743
688	Vera Kaelin, Vivian Villegas, Yi-Fan Chen, Natalie Murphy, Elizabeth Papautsky, Jodi Litfin, Natalie Leland, Varun Maheshwari, Beth McManus, and Mary Khetani. 2022. Effectiveness and scalability of an electronic patient-reported outcome measure and decision support tool for family-centred and participation-focused early intervention: Prospect hybrid type 1 trial protocol. <i>BMJ open</i> , 12(1):e051582.	744
689		745
690		746
691		747
692		748
693		
694		
695		
696	Vera C. Kaelin, Dianna L. Bosak, Vivian C. Villegas, Christine Imms, and Mary A. Khetani. 2021. Participation-focused strategy use among caregivers of children receiving early intervention . <i>American Journal of Occupational Therapy</i> , 75:7501205090.	749
697		750
698		751
699		752
700		753
701	Vera C Kaelin, Andrew D Boyd, Martha M Werler, Natalie Parde, and Mary A Khetani. 2023. Natural language processing to classify caregiver strategies supporting participation among children and youth with craniofacial microsomia and other childhood-onset disabilities. <i>Journal of Healthcare Informatics Research</i> , 7(4):480–500.	754
702		755
703		756
704		757
705		758
706		759
707		760
708	M. A. Khetani, V. Kaelin, S. Rizk, M. Angulo, Z. Salgado, Y. F. Chen, V. Villegas, J. Dooling-Litfin, N. Leland, E. Lerner Papautsky, N. Murphy, B. McManus, and High Value Early Intervention Research Group. 2023. Preliminary effectiveness of an electronic patient-reported outcome measure and decision support tool on early intervention service quality. <i>Developmental Medicine & Child Neurology</i> , 65(S3):5–87.	761
709		762
710		763
711		764
712		765
713		766
714		767
715		768
716	Mary A Khetani, Erin C Albrecht, Jessica M Jarvis, David Pogorzelski, Emmy Cheng, and Karen Choong. 2018a. Determinants of change in home participation among critically ill children. <i>Developmental Medicine & Child Neurology</i> , 60(8):793–800.	769
717		770
718		771
719		772
720		773
721	Mary A Khetani, James E Graham, Patricia L Davies, Mary C Law, and Rune J Simeonsson. 2015. Psychometric properties of the young children’s participation and environment measure. <i>Archives of physical medicine and rehabilitation</i> , 96(2):307–316.	774
722		775
723		776
724		777
725		778
726	Mary A. Khetani and Isabella Lucero. 2024. Classifying children’s participation and its family of related concepts through the participation and environment measure (pem) approach . In Melek and Aslam, editors, <i>ICF Turkey Conference Proceedings</i> , page 47. Atılım University.	779
727		780
728		781
729		782
730		783
731		784
732	Mary A Khetani, Beth M McManus, Erin C Albrecht, Vera C Kaelin, Jodi K Dooling-Litfin, Elizabeth A Scully, and High Value Early Intervention Research Group. 2020. Early intervention service intensity and young children’s home participation. <i>BMC pediatrics</i> , 20(1):330.	785
733		786
734		787
735		788
736		789
737		790
		791
		792
		793
		794
	Mary A Khetani, Beth M McManus, Kristen Arestad, Zachary Richardson, Renee Charlifue-Smith, Cordelia Rosenberg, and Briana Rigau. 2018b. Technology-based functional assessment in early childhood intervention: a pilot study. <i>Pilot and feasibility studies</i> , 4(1):65.	
	Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed prompting: A modular approach for solving complex tasks. <i>arXiv preprint arXiv:2210.02406</i> .	
	Rita Kukafka, Michael E. Bales, Ann Burkhardt, and Carol Friedman. 2006. Human and automated coding of rehabilitation discharge summaries according to the international classification of functioning, disability, and health . <i>Journal of the American Medical Informatics Association</i> , 13(5):508–515.	
	Denis Newman-Griffis, Jonathan Camacho Maldonado, Pei-Shu Ho, Maryanne Sacco, Rafael Jimenez Silva, Julia Porcino, and Leighton Chan. 2021. Linking free text documentation of functioning and disability to the icf with natural language processing. <i>Frontiers in Rehabilitation Sciences</i> , 2:1–17.	
	OpenAI. 2025. Gpt-5.1. https://openai.com/index/gpt-5-1/ . Large language model with enhanced instruction-following and conversational capabilities.	
	Sabrin Rizk, Vera C Kaelin, Julia Gabrielle C Sim, Natalie J Murphy, Beth M McManus, Natalie E Leland, Ashley Stoffel, Lesly James, Kris Barnekow, Elizabeth Lerner Papautsky, and 1 others. 2023. Implementing an electronic patient-reported outcome and decision support tool in early intervention. <i>Applied Clinical Informatics</i> , 14(01):091–107.	
	Saeideh Shahin, Sara Ahmed, Briano DiRezze, and Dana Anaby. 2024. Reliability and validity of the youth and young-adult participation and environment measure (y-pem): An initial evaluation . <i>Physical occupational therapy in pediatrics</i> , 44(2):232–247.	
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	
	Mina Valizadeh, Vera C Kaelin, Mary A Khetani, and Natalie Parde. 2024. Carecorpus: a corpus of real-world solution-focused caregiver strategies for personalized pediatric rehabilitation service design. In <i>Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , Torino, Italy, May 20-25, 2024, pages 2871–2882. ELRA Language Resource Association.	
	Jason Wei and Kai Zou. 2019. Eda: Easy data augmentation techniques for boosting performance on text classification tasks. <i>EMNLP</i> .	

795 World Health Organization. 2007. *International classification of functioning, disability and health: children and youth version: ICF-CY*. World Health Organization.
796
797
798

799 **A Frequent Bigrams by Strategy Type**

800 In Figure 6, we show the most frequent bigrams for
801 each fine-grained strategy type.

802 **B Manual Augmentation Samples**

803 In the following tables, we show manual augmentation
804 samples for each augmentation technique, for
805 each fine-grained caregiver strategy type.



Figure 3: Top 10 most frequent bigrams in Demands of the Activity strategies.

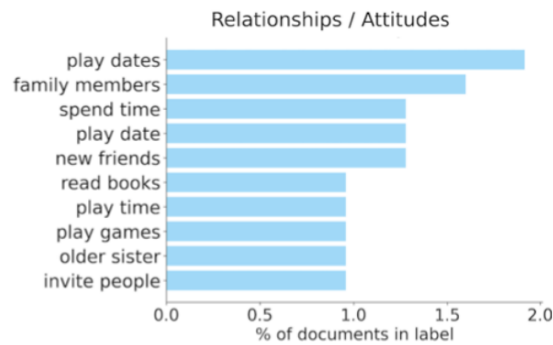


Figure 4: Top 10 most frequent bigrams in Relationships/Attitudes strategies.

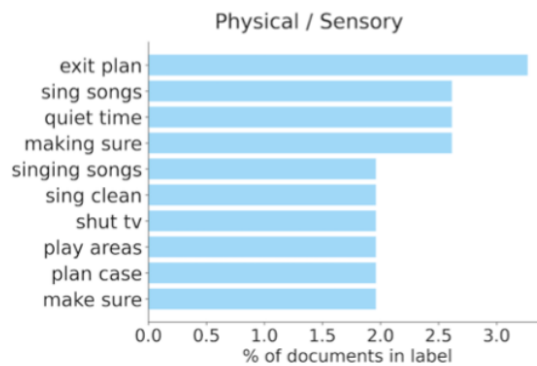


Figure 5: Top 10 most frequent bigrams in Physical/Sensory strategies.

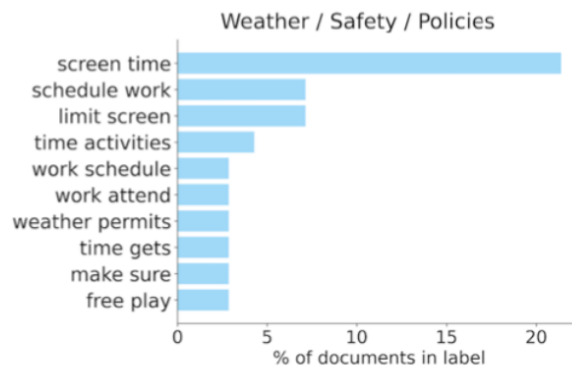


Figure 6: Top 10 most frequent bigrams in Weather/Safety/Policies strategies.

1 = Activity Demands		
Setting	Technique	Example
Home	Original	“Routines and consistency so she knows what to expect”
	D	→ “Routines so she knows what to expect”
	R	→ “Routines to help her predict what will happen”
	A	→ “Set routines so she knows what to expect”
	S	→ “So she knows what to expect, routines and consistency ”
Daycare/Preschool	Original	“Preparing him for what he has to do and what is going to happen”
	D	→ “Helping him get ready for what’s coming”
	R	→ “Preparing him for what to do and what will happen”
	A	→ X
	S	→ “Preparing him for what is going to happen and what he has to do”
Community	D	→ “Short event length so she doesn’t fatigue ”
	R	→ “Reduce the event duration to avoid fatigue”
	A	→ X
	S	→ “So she doesn’t fatigue, keep the length of the event short”

Appendix Table : Manually generated examples for Class 1 (Activity Demands).

2=Relationship/Attitudes		
Setting	Technique	Example
Home	Original	“Grandparents play with him and talk to him”
	D	→ “Grandparents play and talk with him”
	R	→ “Grandparents interact and speak with him”
	A	→ “Have grandparents play and talk with him”
	S	→ “He plays with grandparents and talks with them”
Daycare/Preschool	Original	“Program is geared for interaction with other children”
	D	→ “Program geared for interaction”
	R	→ “Program structured for interaction”
	A	→ “Provide program geared for interaction”
	S	→ X
Community	Original	“Lots of physical support, holding, affection helps”
	D	→ “Physical support, holding, and affection”
	R	→ “Physical assistance, cuddling, and warmth”
	A	→ “Provide physical support, holding, and affection”
	S	→ “Holding, affection, and physical support”

Appendix Table : Manually generated examples for Classes 1-2. (Continued)

3=Physical/Sensory		
Setting	Technique	Example
Home	Original	“Breaks for quiet time away from the hubbub of crowds and noise”
	D	→ “Breaks for quiet time”
	R	→ “Breaks to calm down ”
	A	→ “Have breaks for quiet time”
	S	→ “Breaks away from the hubbub of crowds and noise for quiet time”
Daycare/Preschool	Original	“Keeping him in the front, where he is less likely to be distracted”
	D	→ “Keeping him in the front for less distraction”
	R	→ “Seating him in the front to reduce distraction”
	A	→ X
	S	→ “Where he is less likely to be distracted, keep him in the front”
Community	Original	“Space to remove himself to if he becomes overwhelmed”
	D	→ “Space to remove himself if overwhelmed”
	R	→ “Space to take a break if overwhelmed”
	A	→ “Provide space for him to remove himself if overwhelmed”
	S	→ “If he becomes overwhelmed, remove him to space”

Appendix Table : Manually generated examples for Class 2 (Physical/Sensory).

4=Weather/Safety/Policy		
Setting	Technique	Example
Home	Original	“Started limiting screen time, though harder to do when weather is bad”
	D	→ “Limit screen time”
	R	→ “Restrict screen time”
	A	→ X
	S	→ “Though harder to do when weather is bad, started limiting screen time”
Daycare/Preschool	Original	“Make sure work schedule allows transporting children to activities”
	D	→ “Make sure work schedule allows transporting children”
	R	→ “Ensure work hour enables bringing children to activities”
	A	→ X
	S	→ “To transport children to activities, make sure your work schedule allows it”
Community	Original	“I also take her to parks for the same reason when weather permits”
	D	→ “Take her to parks when weather permits”
	R	→ “Visit parks with her when the weather allows”
	A	→ X
	S	→ “When weather permits, I also take her to parks for the same reason”

Appendix Table : Manually generated examples for Classes 3-4.(Continued)

5=Resources		
Setting	Technique	Example
Home	Original	“PEC system to help with steps for the routine we are trying to do”
	D	→ “PEC system to help with steps for the routine ”
	R	→ “PEC system to support steps in the daily routine”
	A	→ “Have PEC system to help with steps for the routine”
	S	→ “Help with steps for the routine we are trying to do through the PEC system”
Daycare/Preschool	Original	“Specialized seating to allow [name] to sit/work with peers”
	D	→ “Specialized seating to work with peers”
	R	→ “Adaptive seating to work with peers”
	A	→ “Assign specialized seating to work with peers”
	S	→ “Allow [name] to sit/work with peers through specialized seating”
Community	Original	“Pictures of things to buy so she can shop”
	D	→ “Pictures of things to buy”
	R	→ “Images of products to purchase”
	A	→ “Have pictures of things to buy”
	S	→ “So she can shop, have pictures of things to buy”

Appendix Table : Manually generated examples for Class 5.