

PICTOEDUCA: Building a Dataset for Spanish Text-to-Pictogram Generation

Anonymous ACL submission

Abstract

We present PICTOEDUCA, the first large-scale Spanish text-to-pictogram dataset for augmentative and alternative communication (AAC), derived from primary educational materials and grounded in the ARASAAC pictogram repository. The dataset is released with a reproducible pipeline that combines automatic annotation with targeted expert correction, supporting scalable and high-quality corpus construction. We benchmark a rule-based system (ARAWORD) and neural models (T5, LLaMA) under direct text-to-pictogram and two-stage text-to-concept-to-pictogram settings. Results show that the rule-based system remains a strong baseline, while neural models benefit from explicit semantic abstraction, with the two-stage approach improving semantic coherence and reducing ambiguity. We further explore data selection strategies, demonstrating that combining domain similarity with a quality signal yields higher-quality silver data, reduces annotation effort, and improves model performance in low-resource regimes. PICTOEDUCA enables reproducible evaluation and advances Spanish text-to-pictogram research.

1 Introduction

Text-to-pictogram generation is a critical task in Augmentative and Alternative Communication (AAC), where natural language sentences are converted into sequences of pictograms to support comprehension and communication for individuals with cognitive or language impairments (Cabello et al., 2018; Schwab et al., 2020). Pictogram-based communication systems have been widely used in educational and clinical contexts, often leveraging tools such as ARAWORD in Spanish¹ or shared tasks such as ImageCLEFtoPicto in French², enabling systematic evaluation of translation and

sequence prediction models (Norré et al., 2021; Koushik et al., 2024). While these efforts show the potential of data-driven and rule-based approaches, most resources target languages other than Spanish, use small corpora, or focus on AAC tools rather than structured datasets for benchmarking.

A key enabling resource for pictogram-based AAC systems is ARASAAC³, an open and multilingual repository of pictograms that has been widely adopted in educational and clinical settings. ARASAAC pictograms serve as a shared visual vocabulary across languages and underpin a variety of tools and research efforts, including ARAWORD in Spanish and text-to-pictogram systems in French and Dutch (Schwab et al., 2020; Norré et al., 2021). The availability of a standardised pictogram inventory facilitates cross-lingual reuse, semantic alignment, and model comparability, while reducing the cost of symbol design and licensing. However, despite its broad adoption, ARASAAC does not provide large-scale, language-specific parallel corpora, leaving dataset construction and benchmarking as an open challenge, particularly for Spanish.

Despite Spanish being the second most spoken language worldwide, there is currently no large-scale Spanish dataset for text-to-pictogram generation. This absence limits the development of data-driven approaches, prevents reproducible evaluation, and constrains research in educational and accessibility applications for Spanish-speaking populations. In particular, creating high-quality, parallel sentence-pictogram corpora is challenging due to the labor-intensive nature of manual annotation and the subtle semantic nuances required for accurate pictogram representation (Bautista et al., 2017).

To address this gap, we present the first large-scale Spanish text-to-pictogram dataset, derived from Peruvian educational materials. The dataset is complemented by a reproducible construction

¹Available at https://aulaabierta.arasaac.org/araword_inicio.

²Available at <https://www.imageclef.org/2025>.

³Available at <https://arasaac.org/>.

079 pipeline that integrates automatic annotation with
080 targeted expert correction and quality control, pro-
081 viding a generalisable framework for creating com-
082 parable resources in other languages or domains.

083 We benchmark a range of text-to-pictogram gen-
084 eration strategies and models, and investigate a data
085 selection strategy to prioritise high-quality, infor-
086 mative instances. This approach enables the con-
087 struction of compact, high-quality ‘silver’ seed sub-
088 sets, facilitates efficient annotation, and improves
089 model performance. Our results establish baselines
090 and shed light on the impact of selective annota-
091 tion when scaling text-to-pictogram datasets. In
092 summary, we make the following contributions:

- 093 • We release PICTOEDUCA, the first large-
094 scale Spanish text-to-pictogram corpus de-
095 rived from educational materials, supporting
096 reproducible evaluation.
- 097 • We provide a systematic evaluation of rule-
098 based and neural generation strategies.
- 099 • We propose a pipeline for selecting, simplify-
100 ing and annotating sentences that can be easily
101 adapted to other domains and languages.
- 102 • We demonstrate how selective annotation can
103 accelerate corpus creation, identify a high-
104 quality seed dataset, and potentially improve
105 model performance.

106 2 Related Work

107 **Spanish Work** In Spanish, ARAWORD⁴ is a
108 widely used AAC tool displaying ARASAAC pic-
109 tograms alongside text, while AraTraductor uses
110 syntactic and morphological Natural Language Pro-
111 cessing techniques to map sentences to pictograms
112 (Bautista et al., 2017). Despite these resources,
113 large-scale, parallel Spanish text-to-pictogram cor-
114 pora for data-driven models remain unavailable, a
115 gap addressed by our dataset.

116 **Other Languages** In French, the ImageCLEFt-
117 oPicto shared task⁵ provides text- and speech-to-
118 pictogram pairs, evaluated with metrics such as
119 BLEU, METEOR, and PictoER (Macaire et al.,
120 2025). Dutch systems and Arasaac-WN further
121 illustrate cross-lingual and semantic alignment ap-
122 proaches linking text to pictograms (Vandeghinste

⁴Available at https://aulaabierta.arasaac.org/araword_inicio.

⁵Available at <https://www.imageclef.org/2025>.

et al., 2017; Sevens et al., 2015; Schwab et al.,
2020; Norré et al., 2021).

Multilingual Prediction Models Transformer-
based models such as PictoBERT and BERTIm-
bau (Souza et al., 2020) variants demonstrate
the value of sequence prediction for pictogram
generation in multiple languages (Pereira et al.,
2022, 2024). Other work in Bengali and multi-
modal French corpora highlight the importance of
aligned datasets for bidirectional text-to-AAC gen-
eration and speech-to-pictogram tasks (Karmakar
and Sinha, 2024; Macaire et al., 2024).

135 3 PICTOEDUCA

136 3.1 Data Sources

137 The corpus was derived from official Peruvian pri-
138 mary school textbooks published by the Ministry
139 of Education of Peru (MINEDU) and distributed
140 through the Perú Educa platform and the Institu-
141 tional Repository⁶. A total of 18 textbooks were
142 collected from the subjects Communication and
143 Science and Technology, targeting pupils from
144 third to sixth grade of primary education.

145 The selected materials span the years 2020-2023
146 and include both regular and rural educational
147 modalities. All resources are publicly available
148 and were chosen to ensure linguistic consistency,
149 educational relevance, and age-appropriate content.
150 The diversity of years and modalities allows the
151 corpus to capture a broad range of instructional
152 styles and vocabulary commonly encountered in
153 Peruvian primary education.

154 3.2 PICTOEDUCA Building Pipeline

155 Figure 1 illustrates the procedure used to construct
156 the PICTOEDUCA dataset.

157 3.2.1 Preprocessing and Filtering

158 Sentences were automatically extracted from PDF
159 documents using custom web scraping. Following
160 extraction, all sentences were manually validated
161 to ensure grammatical correctness, semantic coher-
162 ence, and relevance to educational contexts. This
163 initial step yielded 27,650 sentences.

164 A subsequent cleaning phase removed dupli-
165 cated entries, sentences containing non-Spanish
166 tokens, and fragments lacking sufficient contextual
167 meaning (e.g. cut segments), reducing the corpus
168 to 21,962 unique sentences.

⁶Available at <https://www.perueduca.pe/#/home/materiales-educativos>.

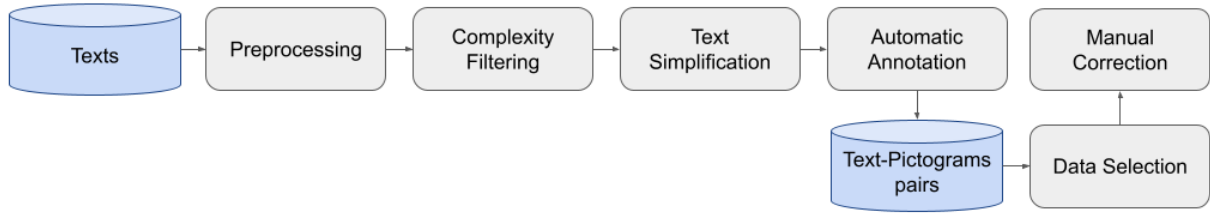


Figure 1: Pipeline for constructing the PICTOEDUCA dataset

Given the target audience of children and individuals with communication difficulties, sentences were categorised according to their linguistic complexity. We adopted the proposal of Vázquez-Rodríguez et al. (2022) for readability assessment in Spanish. Using this framework, sentences were classified into simple and complex categories.

The resulting distribution comprised 17,626 simple sentences and 4,336 complex sentences. Only simple sentences were retained for further processing, as they typically exhibit shorter length, lower lexical density, and more predictable syntactic structures, all of which are beneficial for pictographic translation tasks.

3.2.2 Automated Text Simplification

To further enhance accessibility, the selected “simple” sentences were automatically simplified using GPT-4, which was prompted to enforce easy-to-read principles, a maximum of ten words, and vocabulary suitable for 7–8-year-old children (see Figure 4, Appendix A). This step reduced residual syntactic and lexical complexity while preserving meaning. After cleaning and deduplication, the corpus was reduced to 16,319 simplified sentences.

3.2.3 Annotation Process

Each sentence was annotated with a sequence of pictogram identifiers from the ARASAAC system⁷. ARASAAC provides an open-access pictogram repository, exposed via public APIs⁸, containing metadata for over 13,500 pictograms, including identifiers and associated keywords.

Automatic annotation was performed using the static rule-based system ARAWORD⁹. For each sentence, tokens were lemmatised and matched against the ARASAAC dictionary. When a direct match was unavailable, fallback strategies were applied, such as using the lemma form or preserv-

ing the original word as metadata. Compound pictograms were detected and substituted when applicable. The final output for each instance consisted of the original sentence, the pictogram identifier sequence, and associated metadata. The algorithm is described in Appendix B.

In addition to the automatically annotated corpus, two manually labelled subsets were created for validation and test. First, 1,111 sentences were sampled from the corpus and annotated by expert psychologists specialised in child language, behaviour, and augmentative and alternative communication (AAC). To support this process, a dedicated web-based annotation tool was developed, allowing experts to select pictograms directly from the ARASAAC repository.

Second, 835 sentences were manually extracted from ARASAAC educational materials and annotated following the same linguistic and pedagogical criteria. These two subsets were combined and split evenly into validation and test sets, each containing 973 sentences. All manually annotated instances were excluded from the automatically annotated data, leaving 15,208 sentences for training.

4 Experimental Setup

4.1 Task Definition

The task consists of translating a Spanish sentence into a sequence of pictogram identifiers from the ARASAAC repository. Given an input sentence $x = (w_1, \dots, w_n)$, the goal is to generate an ordered sequence $y = (pict_1, \dots, pict_n)$, where each $pict_n$ corresponds to a pictogram identifier that visually represents part of the sentence meaning. Figure 2 shows an example.

The task is framed as a sequence generation problem and evaluated at the sentence level. Two task formulations are considered:

- Direct text-to-pictogram translation: where the model predicts pictogram identifiers directly from text.

⁷Available at <https://arasaac.org/>.

⁸Available at <https://arasaac.org/developers/api>.

⁹Available at https://aulabierta.arasaac.org/araword_inicio

Text: *Cerrar los ojos y dormir*
(close your eyes and sleep)

Pictogram Identifiers: pict_1014 pict_6479

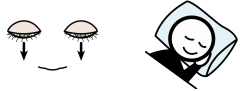
Pictograms: 

Figure 2: Input text, pictogram identifiers and pictograms (images) for the sentence “Cerrar los ojos y dormir”.

- Two-stage text-to-concept-to-pictogram translation: where an intermediate sequence of textual concepts is first generated and then mapped to pictograms.

4.2 Modeling Approaches

To address the task, three different modelling strategies were implemented: a rule-based baseline, a large language model using prompting, and a neural sequence-to-sequence model.

4.2.1 Rule-Based Baseline: ARAWORD

As a baseline, we employ ARAWORD, a rule-based system designed to translate Spanish sentences into ARASAAC pictograms. ARAWORD relies on lemmatisation, dictionary matching, and handcrafted linguistic rules to map words or lemmas to pictogram identifiers.

4.2.2 Large Language Model: LLaMA-8B

The second approach uses LLaMA 3.1–8B-Instruct (Grattafiori et al., 2024) under a two-stage text-to-concept-to-pictogram strategy: the model first predicts key concepts from each sentence, which are then mapped to pictogram identifiers.

The prompt (Figure 5, Appendix C) guides the model to extract visually representable actions, entities, attributes, and relations, following AAC principles (e.g., implicit subjects, generic nouns). The resulting concepts are mapped to ARASAAC pictograms via keyword matching.

To resolve multiple candidate pictograms, a multilingual CLIP model (Reimers and Gurevych, 2019)¹⁰ computes text–image similarity, selecting the pictogram that best aligns with each concept.

4.2.3 Sequence-to-Sequence Model: T5

The third approach employs a neural encoder–decoder architecture based on the Spanish

¹⁰Available at <https://huggingface.co/sentence-transformers/clip-ViT-B-32-multilingual-v1>

T5 Base model (Raffel et al., 2020; Araujo et al., 2024), fine-tuned on the annotated corpus under both the direct text-to-pictogram and two-stage text-to-concept-to-pictogram settings.

In the direct setting, the model generates pictogram identifiers as target tokens, while in the concept-based configuration it predicts textual concepts, which are mapped to pictograms via the same CLIP-based selection used for LLaMA.

In addition, for direct setting, we extend model’s vocabulary to include all ARASAAC identifiers. Training used a batch size of 16, a learning rate of 1×10^{-5} , 40 epochs, and a weight decay of 0.01.

4.3 Automatic Evaluation

System outputs were evaluated against expert-annotated references using standard sequence generation metrics:

- BLEU (Papineni et al., 2002), to measure n-gram overlap between predicted and reference pictogram sequences.
- chrF++ (Popović, 2017), which captures character-level similarity and is more tolerant of near-miss pictogram identifiers.
- n-gram precision (1–4), to analyse performance degradation as sequence length increases.

Table 1 presents the automatic evaluation results for the proposed text-to-pictogram approaches. The rule-based system ARAWORD achieves the best overall performance, obtaining the highest BLEU (0.3741), chrF++ (51.46), and n-gram precision scores. These results confirm ARAWORD as a strong baseline, particularly in scenarios where a well-defined pictographic vocabulary and deterministic mappings are available, highlighting the relevance of rule-based methods in AAC settings.

Neural models benefit from a two-stage pipeline that separates semantic interpretation (Text-to-Concept) from pictogram selection (Concept-to-Pictogram). This is evident for T5, whose BLEU score increases from 0.2977 in the direct setting to 0.3218 in the two-stage configuration, alongside a substantial improvement in chrF++ (42.95 to 48.23). This decomposition appears well suited to pretrained text-based models, as the direct approach requires learning and sequencing pictogram identifiers, which likely demands more data. Aligning pictogram embeddings with pretrained textual representations may help alleviate this issue.

330 Within the two-stage setting, LLaMA-8B
 331 achieves its highest BLEU score (0.3452) under 1-
 332 shot prompting. While additional in-context exam-
 333 ples do not consistently improve BLEU, they yield
 334 moderate gains in chrF++, suggesting improved
 335 selection of semantically related pictograms even
 336 when exact label matches are not achieved (e.g.,
 337 alternative pictograms for “yo” - Figure 3). How-
 338 ever, despite strong unigram precision, LLaMA-8B
 339 shows a marked drop in higher-order n -gram pre-
 340 cision, indicating difficulties in generating coher-
 341 ent longer pictogram sequences. Future work may
 342 explore larger or specialised models and structure-
 343 aware decoding strategies to address this limitation.

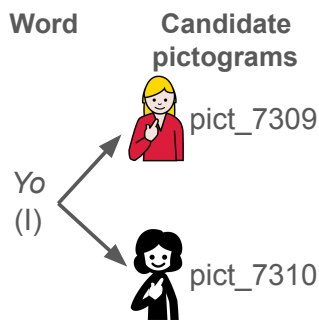


Figure 3: Candidate pictograms and its corresponding pictogram ids for the word “Yo” (I).

344 4.4 Human Evaluation

345 To assess semantic adequacy and interpretability,
 346 a human evaluation was conducted with psycholo-
 347 gists, speech therapists, and special education pro-
 348 fessionals experienced in AAC. In total, 8 experts
 349 participated in the human evaluation.

350 A random sample of 50 test sentences was evalu-
 351 ated. For each instance, evaluators rated the cor-
 352 respondence between the original sentence and the
 353 generated pictogram sequence on a five-point Lik-
 354 ert scale (1 = very poor, 5 = very good). This
 355 evaluation focuses on (semantic) adequacy rather
 356 than strict sequence matching.

357 Table 2 reports the average human evaluation
 358 scores for each system. ARAWORD achieves the
 359 highest mean score (3.04), confirming its role as a
 360 strong baseline. Among neural approaches, T5
 361 benefits substantially from the Text-to-Concept-
 362 to-Pictogram setting, reaching a mean score of
 363 2.79 and clearly outperforming the direct Text-to-
 364 Pictogram variant (2.18). This improvement indi-
 365 cates that introducing an explicit conceptual layer
 366 helps reduce semantic ambiguity in pictogram gen-
 367 eration. LLaMA attains a moderate score (2.42),

368 suggesting reasonable performance, albeit below
 369 that of the Text-to-Concept-to-Pictogram T5.

370 Statistical analysis confirms that these differ-
 371 ences are significant. A Friedman test reveals over-
 372 all variation across systems ($\chi^2 = 20.68$, $p <$
 373 0.05). Post-hoc Wilcoxon tests ($p < 0.05$) with
 374 Holm correction show that ARAWORD signifi-
 375 cantly outperforms both T5 Text-to-Pictogram and
 376 LLaMA, while T5 Text-to-Concept-to-Pictogram
 377 yields a significant improvement over its direct
 378 counterpart. No significant difference is observed
 379 between ARAWORD and the concept-based T5
 380 model, suggesting that semantic abstraction is a
 381 promising direction for narrowing the gap between
 382 neural systems and expert-designed references.

383 It is also worth noting that automatic metrics
 384 show a moderate system-level correlation with
 385 human judgements, capturing broad performance
 386 trends but failing to reflect finer-grained commu-
 387 nicative adequacy. Approaches with higher BLEU
 388 and chrF++ scores, such as ARAWORD and the
 389 Text-to-Concept-to-Pictogram T5, are also pre-
 390 ferred by human evaluators, indicating sensitivity
 391 to gains from semantic abstraction. However, dis-
 392 crepancies arise for neural models: despite strong
 393 unigram precision and competitive BLEU scores,
 394 LLaMA receives lower human ratings, suggest-
 395 ing that n -gram-based metrics overestimate per-
 396 formance by prioritising local overlap over global se-
 397 mantic coherence. Among the automatic measures,
 398 chrF++ aligns more closely with human judgement,
 399 yet none adequately capture acceptable semantic
 400 variants common in AAC, motivating the contin-
 401 ued use of human evaluation and the development
 402 of more semantically grounded metrics.

403 5 Impact of Data Selection on Annotation 404 and Model Performance

405 Annotating sentences with corresponding pic-
 406 togram sequences is labour-intensive, and although
 407 automatic systems such as ARAWORD enable
 408 scalable annotation, they often introduce noise and
 409 fail to capture fine-grained semantic distinctions.
 410 Our initial corpus of 15,208 automatically anno-
 411 tated sentences therefore contained inconsistencies
 412 that made exhaustive manual correction impracti-
 413 cal. In this section, we investigate whether data
 414 selection can mitigate these limitations by guid-
 415 ing annotation effort towards the most informa-
 416 tive instances. Specifically, we examine *whether*
 417 *data selection strategies are effective for identify-*

Setting	Approach	BLEU	Precision-1	Precision-2	Precision-3	Precision-4	chrF++
TEXT-TO-PICTOGRAM	ARAWORD	0.3741	0.6075	0.4721	0.3096	0.2205	51.46
	T5	0.2977	0.5699	0.4107	0.2187	0.1534	42.95
TEXT-TO-CONCEPT-TO-PICTOGRAM	T5	0.3218	0.5766	0.4273	0.2494	<i>0.1745</i>	<i>48.23</i>
	LLaMA-8B 0-SHOT	0.3061	0.6447	0.4563	0.2204	0.1354	41.43
	LLaMA-8B 1-SHOT	<i>0.3452</i>	0.6562	0.4790	<i>0.2581</i>	<i>0.1751</i>	44.60
TEXT-TO-CONCEPT-TO-PICTOGRAM	LLaMA-8B 2-SHOT	0.3195	0.6313	0.4526	0.2330	0.1565	44.53
	LLaMA-8B 3-SHOT	0.3229	0.6278	0.4522	0.2373	0.1614	45.16

Table 1: Automatic evaluation results. **Bold** values indicate the overall highest scores, while *italic* values indicate the highest scores within the Text-to-Concept-to-Pictogram setting. The results reported for T5 and LLaMA are based on only 1 execution.

Approach	Avg. Score
ARAWORD	3.04*
T5 (Text-to-Pictogram)	2.18
T5 (Text-to-Concept-to-Pictogram)	2.79*
LLaMA-8B-1-SHOT	2.42

Table 2: Human evaluation scores. *Indicates no statistically significant difference.

ing high-quality silver data, accelerating the annotation process by prioritising impactful examples, and ultimately contributing to measurable improvements in downstream model performance. To this end, we adopt a pool-based active learning framework to select candidate instances for expert review, with the goal of improving annotation quality while reducing manual effort and constructing a more reliable text-to-pictogram dataset.

5.1 Selection Criteria

We defined two complementary criteria to estimate the informativeness of a sentence for correction:

Domain Similarity via Feature Decay Algorithms (FDA) : To ensure that reviewed instances were representative of the intended application domain (educational texts for children), we applied the Feature Decay Algorithm (FDA) (Biçici and Yuret, 2015). FDA ranks candidate sentences based on their n-gram overlap with a reference set (S_{seed}) drawn from the validation corpus. The score for each candidate sentence s is computed as:

$$score(s, S_{seed}, L) = \frac{\sum_{ngr \in \{s \cap S_{seed}\}} 0.5^{C_L(ngr)}}{length(s)}$$

Where ngr denotes an n-gram, and $C_L(ngr)$ is the count of how often the n-gram has already appeared in the selected set L . This scoring prioritizes sentences that introduce informative, diverse content relative to the validation domain.

Translation Quality Estimation via X-CLIP

(TQE) : As gold-standard pictogram references are unavailable for most of the training data, we propose a proxy metric to estimate translation quality. Pictogram sequences are modelled as visual streams and compared to sentence meaning using X-CLIP (Ni et al., 2022), a multimodal transformer trained for video-text alignment. As X-CLIP primarily operates in English, Spanish sentences are first translated into English using GPT-4. Cosine similarity between the translated text and the pictogram embeddings provides a weak semantic alignment score, with higher values indicating more faithful pictographic translations.

To combine both metrics into a unified selection mechanism, we defined a composite informativeness score for each instance:

$$score(s) = FDA(s).TQE(s)$$

This multiplicative formulation ensures that selected examples are both domain-relevant and pictographically meaningful, filtering out low-impact or low-quality candidates.

5.1.1 Data Selection for better silver data

We evaluate data selection strategies using a Spanish T5 model trained for direct text-to-pictogram generation. The initial training corpus consists of 15,208 automatically annotated sentences. From this pool, we construct subsets of varying sizes (1,000, 3,000, 5,000, and 7,000 sentences) using three sampling strategies:

- **Random**: sentences sampled uniformly at random from the base corpus (baseline);
- **Similarity**: sentences selected based on domain similarity using FDA;
- **Similarity + Quality**: sentences selected using a combination of FDA-based similarity and TQE-based quality estimation.

Table 3 shows the results of the experiments. Across all settings, performance improves with increasing data size, reflecting the benefits of additional training signal. For example, BLEU scores rise from approximately 0.009–0.019 with 1,000 examples to above 0.23 with 7,000 examples, while chrF++ increases from around 20 to over 36.

The selection strategy has a pronounced effect in low-data regimes. Random sampling consistently yields the weakest results, while similarity-based selection provides moderate gains by prioritising domain-relevant examples. The strongest improvements are observed when similarity is combined with quality-based filtering. With only 1,000 examples, Precision-1 increases from 0.332 under Random sampling to 0.481 with Similarity + Quality, demonstrating the value of careful data curation.

As the subset size grows, the performance gap between Similarity and Similarity + Quality narrows, indicating diminishing returns from quality filtering in higher-resource settings. At 7,000 examples, BLEU scores are nearly identical for both strategies (0.228 vs. 0.230), although both outperform Random selection. Overall, these results suggest that quality-aware data selection is particularly effective for constructing high-quality initial training sets, while larger datasets naturally reduce the marginal benefit of additional filtering.

5.1.2 Annotation Efficiency

We estimate annotator effort by computing the word-level edit distance between the model-generated pictogram sequences and their human-corrected versions. This metric reflects the number of insertions, deletions, and substitutions required for a sentence to be usable in practice, with higher values indicating greater annotation effort.

To analyse annotation efficiency under realistic selection conditions, we consider a subset of 1,000 instances selected using two strategies: *Domain Similarity* and *Domain Similarity + Quality*. Since data selection methods typically prioritise examples ranked as most informative, we focus on the *last 250 instances* selected by each strategy. These instances are expected to be the least similar to the validation set; however, under the Similarity + Quality strategy, they are also constrained to satisfy a minimum quality criterion. This setup allows us to assess which strategy better identifies examples suitable for expert review.

Table 4 reports the distribution of the edit distances for both strategies. Results show that incor-

porating quality information consistently reduces annotator effort. The proportion of examples requiring no changes increases to 20% under Similarity + Quality, compared to 14% with Similarity alone, while high-effort cases (edit distance > 0.75) decrease from 12.8% to 9.2%. Mid-range edit distances (0.25–0.75) remain comparable across strategies, indicating that most annotation effort is concentrated on moderately challenging examples.

Overall, these findings suggest that combining domain similarity with a quality-aware signal more effectively filters out low-quality instances, reducing the number of examples that require more manual correction and improving annotation efficiency.

5.1.3 Impact of Manual Correction on Model Performance

To assess how expert intervention affects downstream model performance, we compare models trained on automatically annotated data with those trained on manually corrected instances selected under different data selection strategies.

Using the same subset analysed for annotation efficiency, Table 5 compares the Similarity and Similarity + Quality strategies under both automatic (“auto”) and human-corrected (“corrected”) conditions. Model performance is reported using Precision-1, Precision-2, and chrF++¹¹.

For the Similarity-based strategy, manual correction leads to consistent improvements across all metrics, with Precision-1 increasing from 0.3085 to 0.3793, Precision-2 from 0.1639 to 0.2051, and chrF++ from 19.03 to 19.70. This indicates that similarity-driven selection surfaces a higher proportion of structurally noisy or misaligned annotations, for which human intervention directly improves surface-level alignment and lexical consistency.

In contrast, for the Similarity + Quality strategy, manual correction results in a slight decrease across metrics. Rather than indicating a negative effect of human intervention, this behaviour suggests diminishing returns when a strong quality signal is already present during selection. In this setting, remaining errors tend to be more subtle and semantic in nature, and manual corrections often prioritise communicative adequacy over exact surface overlap. Given the partial nature of correction, this may also introduce minor inconsistencies within an otherwise highly regularised subset, which are not

¹¹BLEU and Precision-3 and Precision-4 are omitted due to near-zero values and limited interpretability.

Subset Size	Strategy	BLEU	Precision-1	Precision-2	Precision-3	Precision-4	chrF++
1000	Random	0.0092	0.3322	0.1768	0.0004	0.0003	20.46
	Similarity	0.0117	0.3034	0.1615	0.0009	0.0004	19.07
	Similarity + Quality	0.0195	0.4813	0.2695	0.0018	0.0006	20.63
3000	Random	0.0502	0.5130	0.2866	0.0082	0.0053	25.80
	Similarity	0.0523	0.5064	0.2823	0.0091	0.0058	26.84
	Similarity + Quality	0.0594	0.5079	0.2842	0.0114	0.0075	27.28
5000	Random	0.1418	0.5189	0.3123	0.0601	0.0415	29.87
	Similarity	0.1617	0.5184	0.3186	0.0761	0.0543	31.01
	Similarity + Quality	0.1631	0.5178	0.3189	0.0776	0.0552	31.21
7000	Random	0.2069	0.5352	0.3472	0.1182	0.0835	34.31
	Similarity	0.2283	0.5421	0.3607	0.1403	0.0989	36.35
	Similarity + Quality	0.2303	0.5416	0.3614	0.1427	0.1008	36.58

Table 3: Automatic evaluation across data selection strategies and subset sizes

Edit Distance Range	Similarity (%)	Similarity + Quality (%)
0.0 (no changes)	14.0	20.0
(0.0, 0.25]	10.8	14.0
(0.25, 0.5]	38.4	37.6
(0.5, 0.75]	24.0	19.2
>0.75	12.8	9.2

Table 4: Distribution of edit distances under different data selection strategies.

well captured by single-reference, n-gram-based metrics such as Precision and chrF++.

Overall, these results highlight that manual correction is most beneficial when applied to data selected by weaker heuristics, where annotation noise is more pronounced. When quality-aware selection already constrains the error space, the marginal gains of partial human correction diminish and may not be reflected by automatic evaluation metrics.

Strategy	Precision-1	Precision-2	chrF++
Similarity (auto)	0.3085	0.1639	19.03
Similarity (corrected)	0.3793	0.2051	19.70
Similarity + Quality (auto)	0.4742	0.2641	20.34
Similarity + Quality (corrected)	0.4091	0.2222	20.31

Table 5: Automatic Evaluation for Similarity and Similarity + Quality selection strategies before and after correct 250/1000 instances.

6 Conclusion and Future Work

We introduced PICTOEDUCA, a large-scale Spanish dataset for text-to-pictogram generation grounded in ARASAAC, along with a reproducible pipeline that combines automatic annotation with targeted expert correction to support scalable dataset construction for AAC.

We benchmarked rule-based and neural approaches under direct and two-stage text-to-

concept-to-pictogram settings. Automatic and human evaluations show that the rule-based system remains a strong baseline, while neural models consistently benefit from explicit semantic abstraction, with the two-stage formulation reducing ambiguity and improving performance.

Finally, we showed that quality-aware data selection improves both annotation efficiency and model performance, particularly in low-resource settings. Combining domain similarity with a proxy quality signal yields higher-quality silver data, reduces annotation effort, and enables stronger models with fewer training examples, with diminishing gains as data size increases.

Future work includes expanding the dataset to additional domains, linguistic variations, and expert annotations to improve generalisation. We also plan to develop evaluation metrics tailored to pictogram sequences that better capture semantic relations and multiple valid realisations, extend the approach to other languages using multilingual resources, and further fine-tune models with human feedback to improve robustness. The dataset provides a reproducible platform for advancing both model development and evaluation in automatic pictogram translation, supporting scalable and socially impactful applications.

Limitations

Despite the contributions of the benchmark and the evaluated models, several limitations remain:

Dataset characteristics : PICTOEDUCA is restricted to Spanish texts from Peru and focuses primarily on educational content, limiting coverage of other domains. In addition, each PICTOEDUCA’s instance has a single reference, which may under-

634	estimate model performance and limit evaluation		
635	variability.		
636	Evaluation metrics : Standard metrics such as		
637	BLEU and chrF++ only partially reflect human		
638	judgement, highlighting the need for more seman-		
639	tically grounded evaluation methods.		
640	Neural model constraints : Models struggle to		
641	learn the “pictogram language” when it contains		
642	unseen identifiers, highlighting the need for strate-		
643	gies to induce embeddings for these tokens. One		
644	potential approach is to initialise pictogram embed-		
645	dings using the embeddings of their definitions or		
646	keywords provided by ARASAAC.		
647	Data selection experiments : Experiments were		
648	conducted on small subsets (250 from 1,000 in-		
649	stances), limiting conclusions on the effect of selec-		
650	tive annotation to better assess model performance		
651	improvements (from subsection 5.1.3). Further-		
652	more, data selection strategies leveraging domain		
653	representativeness (FDA) and quality-based filter-		
654	ing (TQE) proved effective. However, these ap-		
655	proaches rely on proxy measures—particularly X-		
656	CLIP, which is designed for video–text alignment		
657	and may not perfectly capture the semantic fidelity		
658	of pictogram sequences.		
659	Human evaluation scope: Assessments were		
660	limited to a 50-sentence subset evaluated by eight		
661	experts, which may constrain the generalisability of		
662	the results. Larger-scale studies are needed to bet-		
663	ter quantify model improvements. Furthermore, in-		
664	cluding children with special communication needs		
665	could provide stronger evidence of real-world effec-		
666	tiveness, although the current evaluation focused		
667	primarily on adequacy rather than communicative		
668	outcomes.		
669	Information About Use of Artificial		
670	Intelligence (AI) Assistants		
671	AI-based tools were used solely for language pol-		
672	ishing and stylistic refinement of the paper. They		
673	were not used to generate research ideas, experi-		
674	mental results, analyses, datasets, or conclusions.		
675	All content, experimental design, and interpreta-		
676	tions were produced and verified by the authors,		
677	who take full responsibility for the paper.		
678	References		
679	Vladimir Araujo, Maria Mihaela Trusca, Rodrigo Tu-		
680	fiño, and Marie-Francine Moens. 2024. Sequence-		
	to-sequence Spanish pre-trained language models.	681	
	In <i>Proceedings of the 2024 Joint International Con-</i>	682	
	<i>ference on Computational Linguistics, Language</i>	683	
	<i>Resources and Evaluation (LREC-COLING 2024)</i> ,	684	
	pages 14729–14743, Torino, Italia. ELRA and ICCL.	685	
	Susana Bautista, Raquel Hervás, Agustín Hernández-	686	
	Gil, Carlos Martínez-Díaz, Sergio Pascua, and Pablo	687	
	Gervás. 2017. Aratrador: text to pictogram trans-	688	
	lation using natural language processing techniques.	689	
	In <i>Proceedings of the XVIII International Conference</i>	690	
	<i>on Human Computer Interaction</i> , Interacción ’17,	691	
	New York, NY, USA. Association for Computing	692	
	Machinery.	693	
	Ergun Biçici and Deniz Yuret. 2015. Optimizing in-	694	
	stance selection for statistical machine translation	695	
	with feature decay algorithms. <i>IEEE/ACM Transac-</i>	696	
	<i>tions on Audio, Speech, and Language Processing</i> ,	697	
	23(2):339–350.	698	
	L. Cabello, E. Lleida, J. Simon, A. Miguel, and	699	
	A. Ortega. 2018. Text-to-Pictogram Summarization	700	
	for Augmentative and Alternative Communication.	701	
	<i>Procesamiento de Lenguaje Natural</i> , 61:15–22.	702	
	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	703	
	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	704	
	Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-	705	
	ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh	706	
	Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-	707	
	tra, Archie Sravankumar, Artem Korenev, Arthur	708	
	Hinsvark, and 542 others. 2024. The llama 3 herd of	709	
	models. <i>CoRR</i> , abs/2407.21783.	710	
	Piyali Karmakar and Manjira Sinha. 2024. Aiding non-	711	
	verbal communication: A bidirectional language ag-	712	
	nostic framework for automating text to AAC genera-	713	
	tion. In <i>Proceedings of the 21st International Confer-</i>	714	
	<i>ence on Natural Language Processing (ICON)</i> , pages	715	
	324–331, AU-KBC Research Centre, Chennai, India.	716	
	NLP Association of India (NLP AI).	717	
	Avaneesh Koushik, Jithu Morrison, P Mirunalini, and 1	718	
	others. 2024. A transformer based approach for text-	719	
	to-picto generation. In <i>Notebook for the ImageCLEF</i>	720	
	<i>Lab at CLEF 2024</i> , pages 1656–1661.	721	
	Cécile Macaire, Chloé Dion, Jordan Arrigo, Claire	722	
	Lemaire, Emmanuelle Esperança-Rodier, Benjamin	723	
	Lecouteux, and Didier Schwab. 2024. A multimodal	724	
	French corpus of aligned speech, text, and pictogram	725	
	sequences for speech-to-pictogram machine transla-	726	
	tion. In <i>Proceedings of the 2024 Joint International</i>	727	
	<i>Conference on Computational Linguistics, Language</i>	728	
	<i>Resources and Evaluation (LREC-COLING 2024)</i> ,	729	
	pages 839–849, Torino, Italia. ELRA and ICCL.	730	
	Cécile Macaire, Diandra Fabre, Benjamin Lecouteux,	731	
	and Didier Schwab. 2025. Overview of the 2025	732	
	ImageCLEFtoPicto Task -Investigating the Genera-	733	
	tion of Pictogram Sequences from Text and Speech	734	
	Notebook for the ImageCLEF Lab at CLEF 2025. In	735	
	<i>Springer Lecture Notes in Computer Science (LNCS)</i> ,	736	
	Madrid, Spain.	737	

738	Bolin Ni, Houwen Peng, Minghao Chen, Songyang Zhang, Gaofeng Meng, Jianlong Fu, Shiming Xiang, and Haibin Ling. 2022. Expanding language-image pretrained models for general video recognition. In <i>Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IV</i> , page 1–18, Berlin, Heidelberg. Springer-Verlag.	
739		
740		
741		
742		
743		
744		
745		
746	Magali Norré, Vincent Vandeghinste, Pierrette Bouillon, and Thomas François. 2021. Extending a text-to-pictograph system to French and to arasaac. In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)</i> , pages 1050–1059, Held Online. INCOMA Ltd.	
747		
748		
749		
750		
751		
752		
753	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	
754		
755		
756		
757		
758		
759		
760	Jayr Pereira, Rodrigo Nogueira, Cleber Zanchettin, and Robson Fidalgo. 2024. Predictive authoring for brazilian portuguese augmentative and alternative communication. <i>Natural Language Processing</i> , 31(2):535–558.	
761		
762		
763		
764		
765	Jayr Alencar Pereira, David Macêdo, Cleber Zanchettin, Adriano Lorena Inácio de Oliveira, and Robson do Nascimento Fidalgo. 2022. Pictobert: Transformers for next pictogram prediction. <i>Expert Systems with Applications</i> , 202:117231.	
766		
767		
768		
769		
770	Maja Popović. 2017. chrF++: words helping character n-grams. In <i>Proceedings of the Second Conference on Machine Translation</i> , pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.	
771		
772		
773		
774		
775	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i> , 21(1).	
776		
777		
778		
779		
780		
781	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	
782		
783		
784		
785		
786	Didier Schwab, Pauline Trial, Céline Vaschalde, Loïc Vial, Emmanuelle Esperanca-Rodier, and Benjamin Lecouteux. 2020. Providing semantic knowledge to a set of pictograms for people with disabilities: a set of links between WordNet and arasaac: Arasaac-WN. In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 166–171, Marseille, France. European Language Resources Association.	
787		
788		
789		
790		
791		
792		
793		
794		
	Leen Sevens, Vincent Vandeghinste, Ineke Schuurman, and Frank Van Eynde. 2015. Extending a Dutch text-to-pictograph converter to English and Spanish. In <i>Proceedings of SLPAT 2015: 6th Workshop on Speech and Language Processing for Assistive Technologies</i> , pages 110–117, Dresden, Germany. Association for Computational Linguistics.	795 796 797 798 799 800 801
	Fábio Souza, Rodrigo Nogueira, and Roberto Lotufo. 2020. Bertimbau: Pretrained bert models for brazilian portuguese. In <i>Intelligent Systems: 9th Brazilian Conference, BRACIS 2020, Rio Grande, Brazil, October 20–23, 2020, Proceedings, Part I</i> , page 403–417, Berlin, Heidelberg. Springer-Verlag.	802 803 804 805 806 807
	Vincent Vandeghinste, Ineke Schuurman Leen Sevens, and Frank Van Eynde. 2017. Translating text into pictographs. <i>Natural Language Engineering</i> , 23(2):217–244.	808 809 810 811
	Laura Vásquez-Rodríguez, Pedro-Manuel Cuenca-Jiménez, Sergio Morales-Esquivel, and Fernando Alva-Manchego. 2022. A benchmark for neural readability assessment of texts in Spanish. In <i>Proceedings of the Workshop on Text Simplification, Accessibility, and Readability (TSAR-2022)</i> , pages 188–198, Abu Dhabi, United Arab Emirates (Virtual). Association for Computational Linguistics.	812 813 814 815 816 817 818 819
	A Prompt used for text simplification	820
	<div data-bbox="820 1570 1402 1939" data-label="Text"> <pre> Convierte las siguientes oraciones en lectura simple para que puedan ser usadas en pictogramas. Debes seguir las instrucciones descritas a continuación: - La salida debe ser una lista de oraciones más simples, cortas y fáciles de leer. - Las oraciones generadas deben ser comprensibles para un niño de 7 u 8 años. - Las palabras utilizadas en las oraciones deben ser fáciles de entender para un niño de esa edad. - La longitud de cada oración no debe exceder las 10 palabras. </pre> </div>	
	Figure 4: Prompt used for automated text simplification.	

B Rule-based algorithm for automatic text-to-pictogram generation.

Algorithm 1: Pseudocode for automatic text-to-pictogram translation.

Input: Initial corpus of sentences,
ARASAAC pictogram dictionary

Output: Corpus translated into pictogram sequences

```

foreach sentence in corpus do
  Lemmatize the sentence into (word,
  lemma) pairs;
  Initialize empty list of pictograms;
  foreach (word, lemma) pair do
    if word exists in dictionary then
      Add corresponding pictogram to
      list;
      Continue to next pair;
    else
      if lemma exists in dictionary
      then
        Add corresponding
        pictogram to list;
      else
        Add word information to list
        (fallback);
      end
    end
  end
  Check for compound pictograms in
  sentence;
  if found then
    Replace compound pictogram in
    sentence;
  end
  Generate sequence of pictogram
  identifiers from list;
  Store original sentence, translation, and
  metadata;
end

```

Eres un psicólogo o pedagogo experto en educación inclusiva y comunicación aumentativa y alternativa (CAA). Tu tarea consiste en identificar los conceptos clave de una oración ("Input") para facilitar la posterior selección de pictogramas adecuados. Sigue cuidadosamente estas instrucciones:

- Extrae los conceptos clave de la oración: pueden ser acciones, entidades, atributos o relaciones esenciales para el significado general.
- Los conceptos deben ser simples, visuales y adecuados para personas con dificultades de comunicación o aprendizaje.
- Puedes mantener conceptos compuestos si tienen un significado claro como unidad (por ejemplo, jugar fútbol), o separar en elementos individuales si mejora la claridad (jugar, fútbol).
- Si el sujeto está implícito, agrégalo usando el pronombre personal adecuado según la conjugación verbal (por ejemplo, yo, tú, él).
- Sustituye nombres propios por un término genérico correspondiente como niño, niña, hombre o mujer, según el género y contexto.
- En preguntas, conserva como concepto clave el pronombre interrogativo principal (que, quién, dónde, 33 cuál, etc.).
- Ordena los conceptos en una secuencia que represente claramente el significado completo de la oración, manteniendo coherencia gramatical y semántica. La secuencia debe funcionar como una traducción pictográfica.
- Retorna únicamente la lista de conceptos clave, separados por comas. No incluyas explicaciones, títulos ni ningún otro texto adicional.

Figure 5: Prompt used for converting text into a sequence of concepts.

C Prompt used for LLama-based Text-to-Concept Generation