

# Beyond “Using Their Own Words”: Abtractivity Characterization in Summarization

Anonymous ACL submission

## Abstract

In this work, we present an extension of the definition of abstractivity within the scope of the automatic generation of summaries. We propose to join extractivity and abstractivity in a single dimension, where extractivity would be on one side of the dimension and complete abstractivity on the opposite one, but in between, there would be levels of abstractivity. A dataset manually annotated to characterize the level of abstractivity of the summaries and to measure the presence of a set of actions applied to compose the summaries has been built. Using this dataset, a study of the sample distribution in terms of abstractivity, annotator agreement, and correlation between annotations regarding the set of actions is presented. An experimental work with a double objective is carried out; on the one hand, we want to validate our perception that extractivity and complete abstractivity are extreme points of a single dimension with multiple abstractivity levels, and on the other hand, we want to verify if there is an overall correlation between the frequency of the actions used for creating the summary and the level of abstractivity. The results confirm both objectives.

## 1 Introduction

Summarizing is the process of condensing the most relevant information from a document into a single, shorter document, the summary. Initially, the essential information in the article has to be identified. There are two strategies to generate the summary from the selected information. In an extractive approach, the sentences with the selected information are copied directly to the summary. In an abstractive approach, the generated summaries also contain the essential information, but it is “expressed, usually, in the words of the author of the summary” (Nenkova and McKeown, 2011).

Although the first approaches to the problem were extractive, after the emergence of the Trans-

former architecture (Vaswani et al., 2017) and its capabilities, most of the published works have addressed the generation of summaries under abstractive approaches. However, to the best of our knowledge, the characterization of abstractivity within summaries has not been sufficiently studied (Bomasani and Cardie, 2020; Grusky et al., 2018; Kryściński et al., 2018; Jing, 2002). A more detailed and extended characterization of the abstractivity in summaries would help to better understand how abstractive models generate their summaries.

Generally, works related to the evaluation of the level of abstractivity of the generated summaries focus on measuring the appearance of new words in the summaries compared to the summarized documents (Wu et al., 2021; Chen et al., 2021; Fu et al., 2021; Manakul and Gales, 2021; Dou et al., 2021; Zou et al., 2020; Zheng et al., 2020). This strategy conforms to Nenkova and McKeown’s definition of abstractive summaries. However, it is not the only way to produce a summary in “the author’s words”. It is possible to make a summary in which very few new words or expressions are introduced compared to the original document, and yet the main ideas are expressed in a different way (Ahuir et al., 2021). For example, a summary can be written based mainly on the reordering of some segments extracted from the document, with the introduction of very few new elements.

In 2002, Jing conducted a study on the actions that abstraction professionals used to create their abstractive summaries (Jing, 2002). Specifically, he identified the following six actions: sentence reduction, sentence combination, syntactic transformation, lexical paraphrase, generalization/specification, and reordering. Based on the hypothesis that writing an abstractive summary is based on using this set of actions, we can characterize the abstractivity of a text by measuring the presence of each of the six actions.

In this work, we propose an extension of the

042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082

083 definition of abstractivity in the automatic summa- 132  
084 rization area. Although in the literature, the extrac- 133  
085 tive and abstractive approaches have been treated 134  
086 as mutually exclusive (Sun et al., 2024; Varab and 135  
087 Xu, 2023; Liu and Lapata, 2019), we join extrac- 136  
088 tivity and abstractivity in a single dimension, what 137  
089 we call the level of abstractivity. The extractiv- 138  
090 ity would be on one side of the dimension, and 139  
091 the complete abstractivity on the opposite one, but 140  
092 in between, there would be levels of abstractivity. 141  
093 Additionally, we want to characterize the level of 142  
094 abstractivity of a summary and measure the pres- 143  
095 ence of each of the actions identified by Jing. 144

096 The main contributions of this work are: 145

- 097 (i) A dataset has been built that contains 146  
098 document-summary pairs manually annotated 147  
099 in terms of a set of actions (including the 148  
100 Jing’s actions) using a Likert scale: the Char- 149  
101 acterization of the Level of Abstractivity in 150  
102 Summarization (CLAsum) dataset. It is pub- 151  
103 licly available at [https://huggingface.co/](https://huggingface.co/datasets/??) 152  
104 [datasets/??](https://huggingface.co/datasets/??). 153
- 105 (ii) Some analyses have been carried out on the 154  
106 CLAsum dataset: sample distribution in terms 155  
107 of abstractivity, annotator agreement, and cor- 156  
108 relation between annotations in terms of the 157  
109 set of actions. 158
- 110 (iii) To check if there is an overall correlation be- 159  
111 tween the frequency of the actions used for 160  
112 creating the summary and the level of abstrac- 161  
113 tivity, two tasks have been defined: Abstrac- 162  
114 tivity Inducing Features extraction, and Abstrac- 163  
115 tivity Level prediction. Both tasks have been 164  
116 addressed as both classification and regression 165  
117 problems. 166
- 118 (iv) Using the CLAsum dataset, a set of machine 167  
119 learning models have been trained to predict 168  
120 both, the Abstractivity Inducing Features and 169  
121 the Abstractivity Level in summaries. 170
- 122 (v) Using these models, some experimentation is 171  
123 carried out to test how beneficial the inclu- 172  
124 sion of the Abstractivity Inducing Features 173  
125 information is in the Abstractivity Level pre- 174  
126 diction. 175

## 127 2 The CLAsum dataset

### 128 2.1 Sample Gathering

129 With the aim of building an appropriate dataset, 180  
130 we selected the test partitions of two well- 181  
131 known datasets in the summarization area: 182

CNN/DailyMail (See et al., 2017), and XSum 132  
(Narayan et al., 2018). Since we wanted diversity 133  
regarding the abstractivity, we distributed the sam- 134  
ples of both test sets into 5 clusters per source using 135  
the KMeans algorithm and selecting some features 136  
related to abstractivity. Specifically, we used the 137  
following abstractivity indicators: Coverage and 138  
Density (Grusky et al., 2018), Content Reorder- 139  
ing (Ahuir et al., 2021), Abstractivity (p=[2,3]) 140  
(Bommasani and Cardie, 2020), and Novel [2,3,4]- 141  
grams (Kryściński et al., 2018). Therefore, an 8- 142  
component features vector was used to characterize 143  
a sample. 144

From those clusters, we extracted 20 samples 145  
per cluster and source. The final set comprised 146  
100 samples from CNN/DailyMail and another 100 147  
from XSum. To ensure the labeling process, some 148  
restrictions were required: (1) the document should 149  
contain a maximum of 500 words, (2) the summary 150  
should contain a minimum of 38 words, (3) the 151  
proportion of words document/summary should be 152  
at least 2:1. 153

### 154 2.2 Labeling Guideline

Since our main objective was to evaluate how the in- 155  
formation in the document was modified (removing 156  
content, merging sentences, etc.), it is necessary 157  
to ensure that a "summary" is really a summary, 158  
that is, much of its information comes from the 159  
document, although it can provide complementary 160  
information. To detect those supposed summaries 161  
that are not really summaries, we included two pre- 162  
vious questions: question A about the relevance 163  
of the information included in the summary with 164  
respect to the document and question B about the 165  
amount of new information added by the summary. 166  
The abstractivity-related questions were 8, from C 167  
to J. One question about the perception of abstrac- 168  
tivity (question C) and 7 questions for the actions 169  
identified by Jing (Jing, 2002) (from D to J); Gen- 170  
eralization (question H) and Specification actions 171  
(question I) were split to gain information. The 172  
complete guideline can be found in Appendix A. 173

We designed the guideline with a Likert scale. 174  
The number of options would vary from question 175  
to question since some aspects required more gran- 176  
ularity than others. For each question, we added 177  
options until we felt that the possible answers col- 178  
lected enough variability and the annotators would 179  
not be forced to choose one option as a fallback. 180  
The number of options are the following ones: (A) 181  
Relevance of the information in the summary (5 182

options), (B) Amount of novel information within the summary (3 options), (C) Perception of the level of abstractivity (5), (D) Content exclusion (4), (E) Sentence information melting (3), (F) Syntax alteration (3), (G) Synonym usage (3), (H) Generalization usage (4), (I) Specification usage (4), and (J) Content Reordering (3).

Additionally, we included the answer 0 (“Does not apply; it is not a summary.”) for questions C to J (abstractivity-related questions). In that way, the annotators would not be forced to answer the abstractivity-related questions if they do not consider the evaluated text a valid summary.

### 2.3 Labeling Process

The labeling process was conducted by people from our research group, a total of 13 people with a high degree level of studies in Computer Science (9 University professors, 4 PhD students, and 1 Master’s degree student). Additionally, 4 Computer Science degree students collaborated with the labeling process. Thus, 17 volunteers with good English level (but not native speakers) contributed to accomplishing the annotation process.

Since we wanted to build a annotated dataset with more than one set of labels per document-summary pair, we established to obtain 3 different sets of labels per pair, acquiring a total of 600 samples (pair+labels). Also, we pursued to capture the variety of perceptions from groups of people, therefore, we distributed the samples to the annotators, avoiding the coincidence 3-annotators group between document-summary pairs as much as possible.

We provided the annotators with the ANONYM labeling tool<sup>1</sup> (Appendix B) and the guideline. To avoid any bias, no labeling examples or instructions were provided. We only encouraged annotators to agree on whether a document-summary contained an actual summary.

### 2.4 Sample Distribution

Table 1, shows the distribution of pairs that contain an actual summary and which ones do not.

Summary	Not Summary
175	25

Table 1: Distribution of document-summary pairs that contain a summary and which do not contain an actual summary (*not-summary*).

<sup>1</sup><https://github.com/anonym/ur1>

We observe that 12.5% of pairs do not contain an actual summary since they do not contain at least some information extracted from the summarized document. All the *not-summaries* pairs came from the XSum dataset. Since we were studying the abstractivity in summaries, we excluded these 25 pairs from the rest of the study.

Regarding the perception of abstractivity (question C), Fig. 1 shows the distribution per source (where 1 represents the extractivity summarization style and 5 the highest perception of abstractivity).

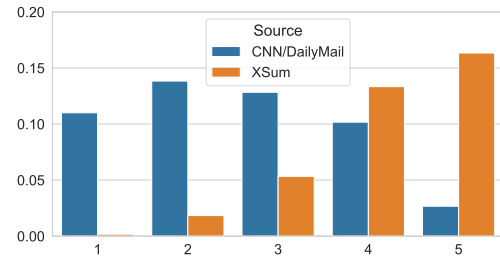


Figure 1: Distribution of answers for question C, regarding the perception of the abstractivity level in the summary.

In Fig. 1, it can be noticed that although the process of selection was the same, the perception of the level of abstractivity from source to source is different. Annotators perceived higher levels of abstractivity in the XSum dataset than in the CNN/DailyMail dataset, which presents more diversity regarding the level of abstractivity.

### 2.5 Annotator Agreement Analysis

The labeling process addresses a complex and subjective task. A total agreement between annotators cannot be expected, then, it would not be advisable to study the agreement in terms of exact matches (binary distance). Therefore, we used the *Relative* distance between two labels.

Eq. (1) shows the definition of this distance.

$$\text{R-Dist}_Q(l_1, l_2) = \frac{|l_1 - l_2|}{M_Q - 1} \quad (1)$$

Given two labels ( $l_1, l_2$ ) for question  $Q$ , *Relative* distance returns the percentage of the absolute distance that separates  $l_1$  from  $l_2$ , in relation to the range between the minimum value (1) and the maximum value that can acquire this question ( $M_Q$ ).

Table 2 shows the average agreement among annotators for each question. We used Cohen’s (Cohen, 1960) and Fleiss’ (Fleiss, 1971) Kappa for the

measurement, with the *Relative* distance (Eq. (1)) as distance function between observations.

Question	Cohen's Kappa	Fleis' Kappa
A	0.94±0.15	0.75±0.21
B	1.00±0.00	0.87±0.22
C	0.92±0.18	0.71±0.19
D	0.92±0.19	0.67±0.23
E	0.96±0.16	0.64±0.34
F	0.90±0.24	0.52±0.30
G	0.90±0.23	0.61±0.28
H	0.86±0.24	0.60±0.22
I	0.86±0.24	0.59±0.22
J	0.89±0.25	0.46±0.32

Table 2: Agreement scores per Question with the *Relative* distance. *Cohen's Kappa* is the pair-wise average score among the three annotators.

It can be observed that the average agreement with Cohen's Kappa is almost perfect. However, when Fleis' Kappa is considered, the agreement strength is reduced to substantial on most of the questions (except B), and moderate for questions F, I, and J. It can be deduced that the annotators' answers do not differ that much for a given question; however, there are slight degree deviations among the three annotations at once (the answers are not unanimous).

We extracted the distances between annotators and questions for each document-summary pair's question to analyze the deviations between annotators. The integer distance was measured between two answers; the distance was computed by counting the number of answers that separated one label from the other. Fig. 2 shows the distribution of integer absolute distance.

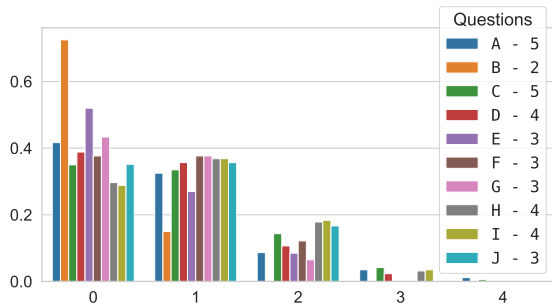


Figure 2: Distribution of answer distances between two annotators on labels for the same document-summary pair.

It can be observed that 30% to 50% of the labels show agreement between annotators, excluding answer B, where the agreement elevates to more than

70% of the cases. However, if we aggregate the annotations with agreement and the ones that are at a distance 1, we cover nearly 80% of the observations in each answer.

With the information extracted from Table 2 and Fig. 2, along with the average Cohen's Kappa between annotators in Appendix C, it can be gathered that the labeling process produced a dataset that captured subjectivity but retained enough agreement to consider the data coherent and valid, from where useful information could be extracted.

## 2.6 Dataset Variants

In the complete dataset, called *Annotators*, for each document-summary pair, there are 3 samples (one per annotator). We also compiled a dataset called *Median*, where the label for a certain question is the median of the 3 corresponding labels.

## 3 Abstractivity-related Questions Correlation Analysis

In this section, we analyze in the CLAsum dataset whether the answers to the questions related to the actions identified by Jing correlate with the perception of the level of abstractivity that the annotators had regarding the viewed summaries.

Fig. 3 presents Pearson's correlation of questions from C to J (abstractivity-related questions) for the *Annotators* dataset. Additionally, we introduce a new column ( $\tilde{x}[D..J]$ ), the median of the 7 aspects (D to J) normalized by the maximum value that can acquire each question.

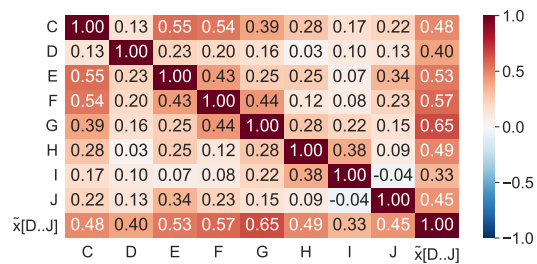


Figure 3: Pearson's correlation between two questions in *Annotators* dataset.  $\tilde{x}[D..J]$  is the normalized median from D to J.

Considering the first column of the matrix (C, perception of abstractivity), two main questions present a moderate correlation with the presence of abstractivity: (E) sentence information melting and (F) syntax alteration, which it is quite clear that it is necessary to create more abstractive summaries.



Using synonyms (G) and generalizations (H) show a low correlation with C, but they still relevant. When we consider the last column, which condenses the perception of the level usage of Jing’s actions, it shows a moderate correlation with C, which means that the perception of abstractivity is related to how frequently those actions were used to compose a summary.

We also studied the correlations using the *Median* dataset, Fig. 4 shows the results.

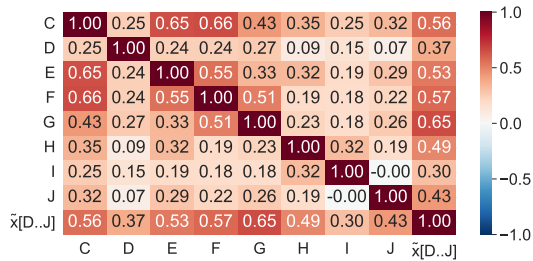


Figure 4: Pearson’s correlation between two questions in *Median* dataset.  $\bar{x}[D..J]$  is the normalized median from D to J.

Questions D and C went up to a high correlation, which is understandable since more abstractive summaries tend to cover more information by joining the information from more sentences, requiring more syntactic changes. The rest of the actions fluctuate between low and moderate correlation. Regarding the last column, the correlation went up closely to the range of high correlation but remained moderate.

All these observations confirm our hypotheses that Jing’s actions are related to the level of abstractivity perception, and, that there exists a single continuous dimension where the two styles, extractive and abstractive summarization, could coexist.

#### 4 Abstractivity Characterization

Based on the conclusions of Section 3, we identify two ways of describing abstractivity in summaries: (1) how often the actions for re-writing and synthesizing the main information from a text have been used to include the information in the summary, and (2) identifying the perception of paraphrasing of the main information of a text included in the summary. This leads us to define two novel tasks regarding the abstractivity in summaries:

**(1) Abstractivity Inducing Features (AIFs) extraction:** Given a document and a summary, the system describes a set of 7 actions in different

grades. (1) Content exclusion [a value from 1 to 4], (2) Sentence information melting [1 to 3], (3) Syntax alteration [1 to 3], (4) Synonym usage [1 to 3], (5) Generalization usage [1 to 4], (6) Specification usage [1 to 4], and (7) Content reordering [1 to 3]. Higher values in a feature indicate a wider presence of a certain action in the summary composition.

**(2) Abstractivity Level (AL) prediction:** Given a document and a summary, the system predicts the level of perception of how much the structure of the document’s main content has been modified to be included in the summary. A value from 1 to 5, where 1 indicates an extractive summarization style and 5 indicates a strong perception that the summary’s author has created it with “their own words”.

The two posed tasks can be approached as ordinal classification or regression problems since both tasks were designed using Likert scales.

### 5 Experimentation

In this section, we detail the experimentation done with both *Median* and *Annotators* datasets. The experimental work has a double objective; on the one hand, we want to validate our perception that extractivity and complete abstractivity are extreme points of a single dimension with multiple abstractivity levels, and on the other hand, we want to verify the role of AIFs in the characterization of these levels of abstractivity.

#### 5.1 Supervised Machine Learning Methods

To tackle the Abstractivity Inducing Features and Level classification/regression tasks, we selected a wide range of classical supervised machine learning methods, all of which were approached with the implementation available in the Scikit-Learn (Pedregosa et al., 2011) Python module.

For classification, the methods that were considered are the following: *Logistic Regression*, *Linear SVM*, *SVM with RFG kernel*, *Random Forest*, and *Multi-Layer Perceptron*. For regression, we used same methods, but *Linear Regression* instead of *Logistic Regression*.

Since some of the methods can not handle more than one feature in the output, we circumvented this handicap by training one model for each feature in the case of AIFs tasks.

#### 5.2 Feature Extraction

Ahvir previously shown (Ahvir et al., 2021) that combining a set of abstractivity-related metrics

(such as *Coverage*, *Density*, *Content Reordering*, *Abtractivity* ( $p=[2,3]$ ), *Novel [2,3,4]-grams*) is useful for abstractivity measurement, and we did not want to include additional variables in the study, then, we followed the same feature extraction as in Section 2.1.

This representation of the document-summary pairs was a straightforward first approach for abstractivity-related tasks. Exploring the impact of other kinds of feature extraction methods, such as incontextual or contextual embeddings, would be out of the scope of the present work.

### 5.3 Evaluation Metrics

We selected a set of metrics for both the classification and regression approaches. For the classification approach, the macro versions of *Precision*, *Recall*, and *F1-score* were used; for the regression approach, we employed *Root Squared Mean Error* (RMSE), *Median Absolute Error* (MdAE).

Additionally, for classification and regression, the *Relative* distance (Eq. (1), Section 2.5) was used in the Abstractivity Level tasks, and the *Minkowski* ( $p=7$ ) distance to measure the distance between the AIFs prediction vector and the reference vector since we want to evaluate the extracted features' cohesion. The *Minkowski* distance between vectors was measured against the normalized AIFs vectors. The normalized AIFs vectors with values from [0, 1] were obtained by dividing each aspect by the maximum value possible for that aspect.

Minimizing *Relative* distance (Abstractivity Level prediction) and *Minkowski* distance (AIFs extraction) will be the main goal to achieve since we want our systems to be as close to the real prediction as possible, and these metrics reflect that need.

### 5.4 System Types Developed

We developed two systems *End-to-End* for each proposed task: one for AIFs classification, another for the regression version of that task, and another two for Abstractivity Level classification and regression. Thus, given a document and a summary, the system first extracts the selected features representing the document-summary pair (Section 5.2), and then performs the classification/regression tasks.

Additionally, we developed a third model type (*AIFs-to-AL*), which receives the document-summary features plus the AIFs as the input and predicts the Abstractivity Level. The model was

trained with the reference AIFs labels as input, along with the document-summary features. This model type should help to analyze how beneficial the inclusion of the AIFs is in the Abstractivity Level prediction.

With the *AIFs-to-AL* models, we created a *Pipeline* for Abstractivity Level prediction. The *Pipeline* receives a document-summary pair, extracts the document-summary features, and with them, the AIFs predictor extracts the corresponding AIFs. Finally, the AIFs are concatenated with the document-summary features and passed to the *AIFs-to-AL* model to obtain the prediction of the Abstractivity Level. The *Pipeline* should help verify the usefulness (for Abstractivity Level prediction task) of a system that considers the AIFs' information compared to a system that does not use them (the *End-to-End* systems).

### 5.5 Training and Evaluation Methodology

With only 175 document-summary pairs to work with, we were facing a low-data situation. For this reason, we trained and evaluated all system configurations 20 times with different partitions, which will show the variability in the performance of each configuration and the conclusions extracted from the results would not be tied to any random aspect of the validation process.

Considering that the distribution of the classes regarding the abstractivity level is not well-balanced, we did not use the K-Fold methodology. Instead, we split the dataset with a different random state (seed) each time. In each partition, 20% of the document-summary pairs were put aside for testing and the rest for training. The partitions were created with the `train_test_split` from Scikit-Learn, setting the seed with an integer number from 0 to 19 and stratified with the C answer (abstractivity level) from *Median* dataset. We verified that all train partitions contain all the possible labels/answers for each question and that 99.4% of the samples were used for testing at least once. Also, it should be mentioned that, in the *Annotators* dataset, all samples that contain the same document-summary pair were placed in the same partition (test or train).

Regarding the sample distribution for training and testing, it should be noted that only one sample per document-summary was available for classification (*Median* dataset). However, for regression (*Annotators* dataset) three samples per document-summary were available, which aimed to capture

the diversity obtained by the labeling process.

For the configuration of each Supervised Machine Learning Method, we bypass modifying the default parameters of the Scikit-Learn implementation (version 1.5.0) to avoid introducing more variables in the study. Only the random state was set to 42 when the method had this feature and increased the max steps to 1 000 000 (a limit that was never reached).

## 6 Systems' Results

This section presents the results obtained by the best system configurations for the two Abtractivity tasks in classification and regression approaches.

All tables follow the same structure. There is a column for each metric, and at the right side of each name, there is an up arrow ( $\uparrow$ ) indicating that a higher value indicates better performance, or a down arrow ( $\downarrow$ ) if lower is better. In each table's numeric cell, the average value and the 95% confidence interval (exponent = lower bound, subscript = upper bound) are shown.

The names of the configurations were shorted for the sake of clearance. The short name are LiR (Linear Regression), LgR (Logistic Regression), LSVM (Linear SVM), Multi-Layer Perceptron (MLP), RnF (Random Forest), and SVM (SVM with RFG kernel). Also, some title names of the columns were shorted: *Mthd* (Method), *M-Dist* (Minkowski distance), *R-Dist* (Relative distance), *Precis* (Precision), and *MdAE* (Median Average Error).

### 6.1 Abtractivity Inducting Features tasks

In this section, the best results for AIFs extraction tasks when we consider Mikowski distance (*M-Dist*) as the reference metric are shown. Extended table results are in Appendix D.

Table 3 shows the best results for the classification task, the *Random Forest* model.

Mthd	M-Dist $\downarrow$	Precis $\uparrow$	Recall $\uparrow$	F1 $\uparrow$
RnF	35.5 <sup>34.7</sup> <sub>36.3</sub>	47.9 <sup>45.0</sup> <sub>50.9</sub>	45.4 <sup>43.6</sup> <sub>47.3</sub>	43.4 <sup>41.5</sup> <sub>45.2</sub>

Table 3: Results of the best model for Abtractivity Inducting Features classification task in *Median* dataset.

Regarding the *M-Dist* average results, we can extract the AIFs predicted vectors average a distance of 36% of the reference vector. The predictions should be considered close enough to be useful, considering that the distance is from a comparison of 7-sized vectors with at least 3 values per feature.

Table 4 shows the best results for AIFs extraction in *Annotators* dataset, the Multi-Layer Perceptron model.

Mthd	M-Dist $\downarrow$	RMSE $\downarrow$	MdAE $\downarrow$
MLP	38.9 <sup>38.4</sup> <sub>39.4</sub>	0.76 <sup>0.75</sup> <sub>0.77</sub>	0.57 <sup>0.56</sup> <sub>0.58</sub>

Table 4: Results of the best model for Abtractivity Inducting Features regression task in *Annotators* dataset.

Relevant results have been achieved for regression. If we consider *RMSE* or *MdAE*, it is noticeable that the model averages less than one level of difference between the predicted feature and the reference one, which indicates that the model can infer a helpful AIFs vector from the abtractivity indicators.

### 6.2 Abtractivity Level tasks

This section presents the best results for each type of system for the Abtractivity Level tasks. The first type is the *End-to-End* (E), the second one is the *Pipeline* (P), and the third model is *AIFs-to-AL* (A). Due to that *AIFs-to-AL* uses the reference AIFs vectors, it can be considered an upper bound of the *Pipeline*.

Table 5 shows the best systems for the *Median* dataset.

Type	Mthd	R-Dist $\downarrow$	Precis $\uparrow$	Recall $\uparrow$	F1 $\uparrow$
E	RnF	10.6 <sup>09.7</sup> <sub>11.5</sub>	64.2 <sup>60.7</sup> <sub>67.6</sub>	60.0 <sup>56.4</sup> <sub>63.6</sub>	60.1 <sup>56.8</sup> <sub>63.5</sub>
P	SVM+RnF	<b>10.2</b> <sup>09.2</sup> <sub>11.2</sub>	<b>64.4</b> <sup>60.3</sup> <sub>68.5</sub>	<b>60.3</b> <sup>56.8</sup> <sub>63.8</sub>	<b>60.3</b> <sup>56.7</sup> <sub>64.0</sub>
A	RnF	9.6 <sup>08.6</sup> <sub>10.6</sub>	68.5 <sup>64.6</sup> <sub>72.4</sub>	63.1 <sup>59.1</sup> <sub>67.0</sub>	63.5 <sup>59.6</sup> <sub>67.5</sub>

Table 5: Results of the best system per system type for Abtractivity Level classification task in *Median* dataset.

Results show that the *Pipeline* system performs slightly better than the *End-to-End* system. This indicates that the AIFs information has positively influenced the performance of the classification task, and it could be improved further if we consider the *AIFs-to-AL* model type scores. However, the *Pipeline* lost 5% of performance due to the cumulated error associated with the AIFs predictor model.

Regarding the regression results, Table 6 shows the results of the best model for each type.

Type	Mthd	R-Dist ↓	RMSE ↓	MdAE ↓
E	ISVM	14.9 <sup>14.2</sup> <sub>15.5</sub>	0.96 <sup>0.93</sup> <sub>1.00</sub>	<b>0.57</b> <sup>0.54</sup> <sub>0.61</sub>
P	SVM+ISVM	<b>14.7</b> <sup>14.1</sup> <sub>15.4</sub>	0.95 <sup>0.92</sup> <sub>0.98</sub>	0.59 <sup>0.56</sup> <sub>0.62</sub>
A	ISVM	13.6 <sup>13.1</sup> <sub>14.1</sub>	0.88 <sup>0.85</sup> <sub>0.91</sub>	0.56 <sup>0.53</sup> <sub>0.59</sub>

Table 6: Restuls of the best system per type for Abstrac- tivity Level regression task in *Annotators* dataset.

In regression, we observe a similar trend as in classification. The information from the AIFs was beneficial for abstractivity level prediction. However, the impact of the AIFs predictor was more noticeable than in classification if we consider the difference in *Pipeline* performance and the *AIFs-to-AL* model.

Since we have observed that AIFs information benefits Abstractivity Level prediction, we compare the confusion matrices (CM) for the 20 runs and 35 test samples per run (700 samples). The number in the y-axis is the reference label, and the one in the x-axis is the predicted. Numbers in cells indicate the number of samples in each combination. Colors indicate the percentage of samples in the combination regarding the total of samples in each row (real label).

Fig. 5 shows the CMs of *End-to-End* model (a), and *Pipeline* model (b).

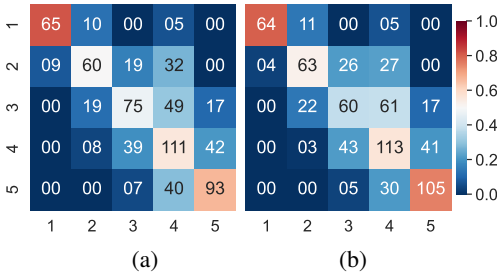


Figure 5: Confusion matrices of *End-to-End* (a) and *Pipeline* (b) in the *Median* dataset for Abstractive Level classification.

Firstly, it is noticed that both systems could correctly classify a sensible number of samples on each level, confirming that models can capture humans' Abstractivity Level perception in summaries. When we compare both systems, (a) and (b), we notice that levels of abstractivity 2, 4, and 5 have increased the number of correct samples (diagonal). Also, the number of samples mislabeled by more than one level has been reduced in levels 4 and 5. However, in level 3, the (b) model has reduced the number of correct samples, misleading level 3 with level 4.

When we compare the CMs of *End-to-End* and *AIFs-to-AL* models in Fig. 6, we obtain similar conclusions.

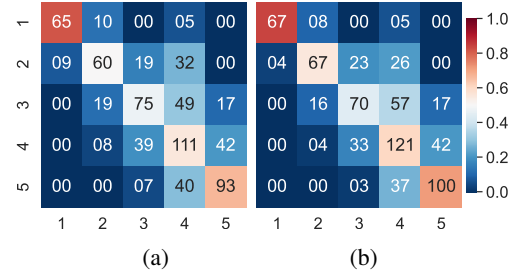


Figure 6: Confusion matrices of *End-to-End* (a) and *AIFs-to-AL* (b) in the *Median* dataset for Abstractive Level classification.

We observe that model (b) increased the number of hits in level 1, in addition to levels 2, 4, and 5. Also, the number of no-hits further than 1 level has been reduced even more than in the *Pipeline* (Fig. 5.b). Finally, the impact on level 3 was less prominent than in the *Pipeline*, but still has lower performance than *End-to-End* for this level.

Generally speaking, we can conclude that AIFs have provided useful information to improve Abstractivity Level prediction, which indicates that measuring these aspects gives additional details about how the summary was composed. Regardless of whether AIFs information was used or not, Fig. 5 shows that models could infer the Abstractivity Level for many samples, supporting the idea of a single continuous dimension where extractive and abstractive summarization coexist.

## 7 Conclusions

In this work, we have presented and made available to the scientific community the CLAsum dataset. This is a hand-annotated dataset that allows characterizing the complexity of the process of summarizing a document by measuring the Abstractivity Level and seven Abstractivity Inducing Features.

The results from the study of the dataset and the experimental work show how the Abstractivity Level and AIFs are related and how AIFs are useful when measuring the level of abstractivity of a summary. Our study places extractivity and complete abstractivity as the extreme points of a single dimension with multiple levels.

## Limitations

**Distribution Representativeness.** The chosen datasets for the study were focused on the news



648 field. Additionally, the selection of document-  
649 summary pairs has been restricted in the number  
650 of words (especially the document side). These  
651 restrictions have reduced the variety of topics that  
652 appeared during the annotation. Consequently, the  
653 distribution of the annotated samples could not be  
654 fully representative of other fields.

655 **Annotators diversity.** Even though 17 anno-  
656 tators from different degrees of studies and expe-  
657 rience have been involved, all are from the field  
658 of Computer Science. Therefore, opinion diver-  
659 sity would be reduced in other fields of expertise  
660 and/or education levels. Additionally, all the an-  
661 notators were not native English speakers. Even  
662 though annotators had high English reading skills,  
663 the fact of not being native speakers, in some spe-  
664 cific situations, a little lost in comprehension of  
665 some particular details of the texts could appear.

666 **Biases.** During all the phases of the annotation  
667 process (guideline design, annotation process, and  
668 data gathering), one of the highest priorities was to  
669 avoid any influence on the outcome. However, we  
670 are mindful that there would always be a chance,  
671 even tiny, that unconscious actions or word selec-  
672 tion could introduce biases in the outcome. In this  
673 regard, we believe that our work produced signifi-  
674 cantly unbiased data that the community could take  
675 as a foundation for future work.

676 **Model Design Soundness.** This work tested a  
677 set of configurations in the most straightforward  
678 possible way to reduce the number of study vari-  
679 ables and presented a basic baseline for future  
680 works. Using a set of abstractivity indicators to  
681 represent the document-summary pairs for the two  
682 proposed abstractivity-related tasks was a direct ap-  
683 proach. In this regard, using them all at once would  
684 not guarantee the best outcome possible with those  
685 indicators since there might be duplicated informa-  
686 tion. Therefore, a correlation study between them  
687 and the abstractivity level would help to reduce  
688 the dimensionality of the features, which could  
689 increase the performance in the tasks. The same  
690 would apply to the Abstractivity Inducting Features  
691 (AIFs), when they are joined to the rest of the indi-  
692 cators and used to predict the level of abstractivity  
693 in the *Pipeline* systems. Additionally, the selection  
694 of supervised machine learning methods was made  
695 without any specific criteria to guarantee the best  
696 outcome; the selection was broadly made to capture  
697 different machine learning method approximations.

## Ethical Statement

**Biases.** The datasets from where the document-  
summary pairs were extracted could present biases  
regarding the vocabulary used or how certain top-  
ics were treated. We did not analyze the included  
samples in this regard; we took them randomly, as  
they were published in the original datasets. Re-  
garding to the findings and statistics presented, they  
are related to annotation complexity and tied to the  
specific group of annotators involved in the process  
and should, therefore, be considered as approxi-  
mate.

**Intended Use and Potential Misuse.** In relation  
to the dataset created in this work, it was created to  
provide the community with data for working and  
expanding the concept of abstractivity in summa-  
rization and new ways to characterize the aspect  
in summaries. Any different analyses or extrapola-  
tions extracted from that data would not be linked  
to the subject of this work and could raise ethical  
considerations.

## References

- Vicent Ahuir, Lluís-F. Hurtado, José Ángel González,  
and Encarna Segarra. 2021. *NASca and NASes: Two  
monolingual pre-trained models for abstractive sum-  
marization in catalan and spanish*. *Applied Sciences*,  
11(21).
- Rishi Bommasani and Claire Cardie. 2020. *Intrinsic  
evaluation of summarization datasets*. In *Proceed-  
ings of the 2020 Conference on Empirical Methods  
in Natural Language Processing (EMNLP)*, pages  
8075–8096, Online. Association for Computational  
Linguistics.
- Xiuying Chen, Hind Alamro, Mingzhe Li, Shen Gao, Xi-  
angliang Zhang, Dongyan Zhao, and Rui Yan. 2021.  
*Capturing relations between scientific papers: An  
abstractive model for related work section generation*.  
In *Proceedings of the 59th Annual Meeting of the  
Association for Computational Linguistics and the  
11th International Joint Conference on Natural Lan-  
guage Processing (Volume 1: Long Papers)*, pages  
6068–6077, Online. Association for Computational  
Linguistics.
- Jacob Cohen. 1960. *A coefficient of agreement for  
nominal scales*. *Educational and Psychological Mea-  
surement*, 20(1):37–46.
- Zi-Yi Dou, Pengfei Liu, Hiroaki Hayashi, Zhengbao  
Jiang, and Graham Neubig. 2021. *GSum: A gen-  
eral framework for guided neural abstractive summa-  
rization*. In *Proceedings of the 2021 Conference of  
the North American Chapter of the Association for*

749			
750		<i>Computational Linguistics: Human Language Technologies</i> , pages 4830–4842, Online. Association for Computational Linguistics.	
751			
752	Joseph L. Fleiss. 1971. <a href="#">Measuring nominal scale agreement among many raters</a> . <i>Psychological Bulletin</i> , 76(5):378–382.		
753			
754			
755	Xiyan Fu, Yating Zhang, Tianyi Wang, Xiaozhong Liu, Changlong Sun, and Zhenglu Yang. 2021. <a href="#">RepSum: Unsupervised dialogue summarization based on replacement strategy</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6042–6051, Online. Association for Computational Linguistics.		
756			
757			
758			
759			
760			
761			
762			
763			
764	Max Grusky, Mor Naaman, and Yoav Artzi. 2018. <a href="#">Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies</a> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 708–719, New Orleans, Louisiana. Association for Computational Linguistics.		
765			
766			
767			
768			
769			
770			
771			
772	Hongyan Jing. 2002. <a href="#">Using hidden markov modeling to decompose human-written summaries</a> . <i>Comput. Linguist.</i> , 28(4):527–543.		
773			
774			
775	Wojciech Kryściński, Romain Paulus, Caiming Xiong, and Richard Socher. 2018. <a href="#">Improving abstraction in text summarization</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1808–1817, Brussels, Belgium. Association for Computational Linguistics.		
776			
777			
778			
779			
780			
781	Yang Liu and Mirella Lapata. 2019. <a href="#">Text summarization with pretrained encoders</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3730–3740, Hong Kong, China. Association for Computational Linguistics.		
782			
783			
784			
785			
786			
787			
788	Potsawee Manakul and Mark Gales. 2021. <a href="#">Long-span summarization via local attention and content selection</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6026–6041, Online. Association for Computational Linguistics.		
789			
790			
791			
792			
793			
794			
795			
796	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. <a href="#">Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807. Association for Computational Linguistics.		
797			
798			
799			
800			
801			
802			
803	Ani Nenkova and Kathleen McKeown. 2011. <a href="#">Automatic summarization</a> . <i>Foundations and Trends® in Information Retrieval</i> , 5(2–3):103–233.		
804			
805			
	F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. <a href="#">Scikit-learn: Machine learning in Python</a> . <i>Journal of Machine Learning Research</i> , 12:2825–2830.		806 807 808 809 810 811 812
	Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. <a href="#">Get to the point: Summarization with pointer-generator networks</a> . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.		813 814 815 816 817 818 819
	Roman Sirokov. 2024. <a href="#">pywebview: A lightweight cross-platform library to create web-based desktop guis</a> .		820 821
	Weisong Sun, Chunrong Fang, Yuchen Chen, Qunjun Zhang, Guan hong Tao, Yudu You, Tingxu Han, Yifei Ge, Yuling Hu, Bin Luo, and Zhenyu Chen. 2024. <a href="#">An extractive-and-abstractive framework for source code summarization</a> . <i>ACM Trans. Softw. Eng. Methodol.</i> , 33(3).		822 823 824 825 826 827
	Daniel Varab and Yumo Xu. 2023. <a href="#">Abstractive summarizers are excellent extractive summarizers</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 330–339, Toronto, Canada. Association for Computational Linguistics.		828 829 830 831 832 833
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all you need</a> . In <i>Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17</i> , page 6000–6010, Red Hook, NY, USA. Curran Associates Inc.		834 835 836 837 838 839 840
	Wenhao Wu, Wei Li, Xinyan Xiao, Jiachen Liu, Ziqiang Cao, Sujian Li, Hua Wu, and Haifeng Wang. 2021. <a href="#">BASS: Boosting abstractive summarization with unified semantic graph</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6052–6067, Online. Association for Computational Linguistics.		841 842 843 844 845 846 847 848 849
	Chujie Zheng, Kunpeng Zhang, Harry Jiannan Wang, Ling Fan, and Zhe Wang. 2020. <a href="#">Topic-guided abstractive text summarization: a joint learning approach</a> . <i>arXiv preprint arXiv:2010.10323</i> .		850 851 852 853
	Yanyan Zou, Xingxing Zhang, Wei Lu, Furu Wei, and Ming Zhou. 2020. <a href="#">Pre-training for abstractive document summarization by reinstating source text</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3646–3660, Online. Association for Computational Linguistics.		854 855 856 857 858 859 860

861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906

## A Labeling Guideline

In this section, the completed guideline that was used in the labeling process is presented.

**Given a newspaper article and a summary in the left side (Document/Summary tab), answer 10 questions/statements regarding the content of the article and the summary and/or the way the summary was created. The possible answers are detailed on the left side (Questions tab).**

A) *The summary provides the most relevant information about the article, and the article extends it with additional details:*

- 0: Strongly Disagree.
- 1: Disagree.
- 2: Undecided.
- 3: Agree.
- 4: Strongly Agree.

B) *Regarding information contained in the summary:*

- 1: All the information in the summary can be found in the article (not necessarily in the exact words).
- 2: Almost all the information in the summary can be found in the article, but adds some additional information.
- 3: I can not consider the given summary a truly abstract. All the information provided in the summary, it is additional and can not be extracted or inferred from the article.

C) *What is your perception about how the author of the summary wrote it?:*

- 0: Does not apply; it is not a summary.
- 1: They rely entirely on the article. It is as if I was reading complete sentences highlighted in the article.
- 2: They rely heavily on the article to write the summary. It only presents slight changes in form and/or order concerning the article.
- 3: They mainly rely on the article to write the summary. Segments of the summary can be identified in the article. Still, the author alters the article's text in form and/or order.

- 4: They weakly rely on the article to write the summary and alter a lot of the article in form and/or order. 907  
908
- 5: Overall, they do not rely on the article to write the summary; instead, they explain the main ideas of the article in their own words. 909  
910  
911  
912  
913

D) *How does the author handle non-relevant information in the article?:* 914  
915

- 0: Does not apply; it is not a summary. 916
- 1: They discard complete sentences. No segments or words of a sentence are discarded. 917  
918  
919
- 2: They focus on mainly discarding complete sentences. Segments or words of sentence discarding is also present, but it is less often than complete sentences discarding. 920  
921  
922  
923  
924
- 3: They focus mainly on discarding text segments within the sentences of the article. The complete sentence discarding is absent, or it is noticeably less frequent than segment. 925  
926  
927  
928  
929
- 4: All information is considered relevant; they manage to cover all the information in the article and substantially reduce its length. discarding. 930  
931  
932  
933

E) *For the creation of the summary, part of the information selected from the sentences of the article is combined to form the sentences of the summary:* 934  
935  
936  
937

- 0: Does not apply; it is not a summary. 938
- 1: No sentences from the article are combined. Each sentence in the summary corresponds to the information contained by a sentence in the article. 939  
940  
941  
942
- 2: Some sentences in the summary are created by combining the information contained by certain sentences from the article. 943  
944  
945  
946
- 3: Most of the sentences of the summary are created by combining information from some sentences of the article. discarding. 947  
948  
949

F) *Sentences in the article that contain the information reflected in the summary have been syntactically altered for inclusion in the summary:* 950  
951  
952  
953

- 0: Does not apply; it is not a summary. 954

955	1: No syntactic alterations exist to create the summary.	3: At most, half of the information susceptible to specification was specified; the rest was not specified.	1002
956			1003
957	2: There are some syntactic alterations to create the summary.	4: More than half of the information susceptible to specification was specified.	1004
958			1005
959	3: There are many syntactic alterations to create the summary.		1006
960			
961	<b>G) When including sentences or segments of the article in the summary, the author replaces words or expressions with semantically equivalent ones:</b>	<b>J) The author of the summary rearranges the chosen information. For example, if facts A-B-C appear in the article, the author refers to them in the following order B-A-C in the summary:</b>	1007
962			1008
963			1009
964			1010
965	0: Does not apply; it is not a summary.	0: Does not apply; it is not a summary.	1012
966	1: Never.	1: Never.	1013
967	2: Sometimes.	2: On one occasion.	1014
968	3: Quite often.	3: On several occasions.	1015
969	<b>H) The summary includes generalizations of information extracted from the article. A generalization is describing one or more concepts using a less specific word (e.g., “Matthew and Amanda reappear in the new sequel of the acclaimed fiction movies of galactic adventures series” in the summary “Matthew and Amanda” could be grouped as “The main actors ...”):</b>	<b>B ANONYM: Anonymized App Name</b>	1016
970			
971		Fig. 7 presents the labeling application developed for the labeling process called ANONYM (Anonymized App Name). The application was developed with Python 3 and PyWebview (Sirokov, 2024), a framework for developing GUI applications with HTML and CSS. The application would be capable of handling different labeling text tasks by just developing an HTML web page for the task needs (supports HTML with CSS and JavaScript). ANONYM is available as a Python module.	1017
972			1018
973			1019
974			1020
975			1021
976			1022
977			1023
978	0: Does not apply; it is not a summary.		1024
979	1: No information can be considered susceptible to generalization without a significant loss of information.		1025
980			1026
981			1027
982	2: No information susceptible to generalization was generalized.		1028
983			1029
984	3: Less than half of the information susceptible to generalization was generalized; the rest was not generalized.		1030
985			1031
986			1032
987	4: More than half of the information susceptible to generalization was generalized.		1033
988			1034
989	<b>I) The summary includes specifications of information extracted from the article. A specification would be to use expressions or words that make the information more specific (e.g., “The race driver has won his ninth F1 World Championship Grand Prix” in the summary “The race driver” could be detailed as “The F1 driver ...”):</b>		1035
990			1036
991			1037
992			
993			
994			
995			
996			
997	0: Does not apply; it is not a summary.		
998	1: No information can be considered susceptible to specification.		
999			
1000	2: No information susceptible to specification was specified.		
1001			

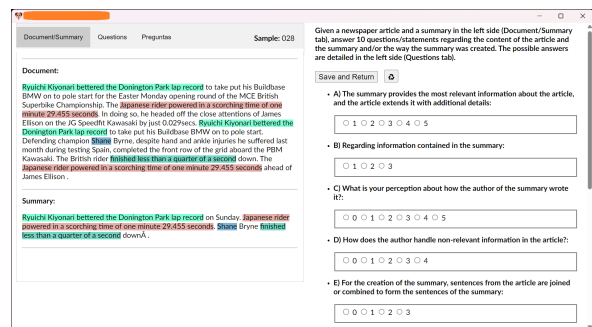


Figure 7: Labeling window of a sample in the ANONYM application.



## C Average Pair-wise Annotator Agreement

Fig. 8 shows the average Cohen’s Kappa agreement between two given annotators. White spaces are combinations that did not occur in the labeling process. The agreement is measured with the *Relative distance* (Eq. (1), Section 2.5) between two annotators and the 10 questions at once.

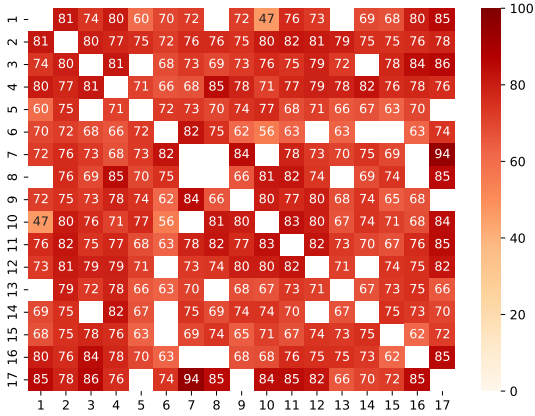


Figure 8: Average of Cohen’s Kappa pair-wise agreement score (*Relative distance*).

## D Extended Results for Abtractivity Inducing Features tasks

Table 7 details the results obtained by all models for the Abtractivity Inducing Features classification task with *Median* dataset. System configurations are sorted in ascending order by the *M-Dist* column.

Mthd	M-Dist ↓	Precis ↑	Recall ↑	F1 ↑
RnF	<b>35.5</b> <sub>34.7</sub> 36.3	47.9 <sub>45.0</sub> 50.9	45.4 <sub>43.6</sub> 47.3	43.4 <sub>41.5</sub> 45.2
SVM	38.4 <sub>36.9</sub> 39.8	41.6 <sub>38.4</sub> 44.7	42.1 <sub>40.7</sub> 43.5	37.7 <sub>36.1</sub> 39.3
LgR	39.6 <sub>38.6</sub> 40.6	48.2 <sub>45.2</sub> 51.3	45.0 <sub>43.4</sub> 46.6	42.2 <sub>40.6</sub> 43.9
ISVM	40.7 <sub>39.9</sub> 41.6	46.6 <sub>43.8</sub> 49.3	43.7 <sub>42.3</sub> 45.0	40.7 <sub>39.4</sub> 42.1
MLP	41.9 <sub>40.7</sub> 43.0	<b>48.5</b> <sub>46.9</sub> 50.2	<b>45.8</b> <sub>44.0</sub> 47.5	<b>44.9</b> <sub>43.3</sub> 46.4

Table 7: Results of models for Abtractivity Inducing Features classification task in *Median* dataset.

Table 8 shows the results obtained by all models for the AIFs regression task with *Annotators* dataset.

Mthd	M-Dist ↓	RMSE ↓	MdAE ↓
MLP	<b>38.9</b> <sub>38.4</sub> 39.4	<b>0.76</b> <sub>0.75</sub> 0.77	<b>0.57</b> <sub>0.56</sub> 0.58
LiR	39.0 <sub>38.6</sub> 39.5	0.77 <sub>0.75</sub> 0.78	0.59 <sub>0.58</sub> 0.60
ISVM	40.6 <sub>40.1</sub> 41.0	0.78 <sub>0.77</sub> 0.79	0.59 <sub>0.57</sub> 0.60
RnF	40.8 <sub>40.3</sub> 41.3	0.79 <sub>0.78</sub> 0.80	0.58 <sub>0.56</sub> 0.60
SVM	40.9 <sub>40.3</sub> 41.5	0.78 <sub>0.77</sub> 0.79	0.55 <sub>0.53</sub> 0.56

Table 8: Results of models for AIFs regression task in *Annotators* dataset.

## E Extended Results for Abtractivity Level tasks

Table 9 details the results for Abtractivity Level classification in *Median* dataset obtained by all configurations for *AIFs-to-AL* (A), *End-to-End* (E) systems, and top-10 configurations for *Pipeline* (P) systems. Numbers in bold are the best average values in their columns, excluding type A type systems since they are not models that can work with the document-summary text, they need the AIFs information.

Type	Mthd	R-Dist ↓	Precis ↑	Recall ↑	F1 ↑
A	RnF	9.6 <sub>08.6</sub> 10.6	68.5 <sub>64.6</sub> 72.4	63.1 <sub>59.1</sub> 67.0	63.5 <sub>59.6</sub> 67.5
P	SVM+RnF	<b>10.2</b> <sub>09.2</sub> 11.2	64.4 <sub>60.3</sub> 68.5	60.3 <sub>56.8</sub> 63.8	60.3 <sub>56.7</sub> 64.0
P	RnF+LgR	10.3 <sub>09.5</sub> 11.2	64.3 <sub>60.9</sub> 67.7	57.7 <sub>54.3</sub> 61.1	58.2 <sub>54.9</sub> 61.5
P	ISVM+RnF	10.4 <sub>09.4</sub> 11.4	<b>64.7</b> <sub>60.6</sub> 68.7	<b>60.8</b> <sub>57.2</sub> 64.4	<b>60.6</b> <sub>56.9</sub> 64.3
P	MLP+LgR	10.4 <sub>09.4</sub> 11.4	62.6 <sub>57.9</sub> 67.2	57.9 <sub>53.5</sub> 62.3	58.0 <sub>53.8</sub> 62.1
P	MLP+RnF	10.5 <sub>09.5</sub> 11.4	64.4 <sub>60.4</sub> 68.3	60.0 <sub>56.7</sub> 63.4	60.2 <sub>56.8</sub> 63.6
P	LgR+RnF	10.5 <sub>09.5</sub> 11.5	63.3 <sub>59.5</sub> 67.0	60.1 <sub>56.7</sub> 63.5	59.8 <sub>56.3</sub> 63.3
E	RnF	10.6 <sub>09.7</sub> 11.5	64.2 <sub>60.7</sub> 67.6	60.0 <sub>56.4</sub> 63.6	60.1 <sub>56.8</sub> 63.5
E	LgR	10.6 <sub>09.6</sub> 11.6	62.8 <sub>58.4</sub> 67.3	57.5 <sub>53.9</sub> 61.1	56.9 <sub>53.4</sub> 60.4
P	RnF+RnF	10.7 <sub>09.8</sub> 11.5	64.3 <sub>61.0</sub> 67.7	59.4 <sub>55.9</sub> 62.9	59.8 <sub>56.5</sub> 63.1
P	LgR+LgR	11.0 <sub>10.0</sub> 12.0	63.1 <sub>59.0</sub> 67.2	58.2 <sub>54.6</sub> 61.8	57.7 <sub>54.1</sub> 61.3
P	ISVM+LgR	11.0 <sub>10.0</sub> 12.1	62.6 <sub>58.2</sub> 67.1	57.5 <sub>53.9</sub> 61.2	57.3 <sub>53.5</sub> 61.1
A	LgR	11.1 <sub>10.2</sub> 11.9	58.5 <sub>54.7</sub> 62.3	54.5 <sub>51.5</sub> 57.6	54.9 <sub>51.7</sub> 58.1
P	RnF+ISVM	11.1 <sub>10.1</sub> 12.1	61.1 <sub>56.7</sub> 65.5	56.4 <sub>52.5</sub> 60.3	55.8 <sub>51.9</sub> 59.6
E	ISVM	11.3 <sub>10.4</sub> 12.1	61.5 <sub>56.6</sub> 66.4	55.8 <sub>52.6</sub> 59.0	54.4 <sub>50.9</sub> 57.8
A	ISVM	11.3 <sub>10.2</sub> 12.5	56.8 <sub>52.4</sub> 61.3	54.3 <sub>50.0</sub> 58.5	53.9 <sub>49.6</sub> 58.1
A	MLP	12.0 <sub>10.7</sub> 13.3	58.3 <sub>53.6</sub> 62.9	54.3 <sub>49.8</sub> 58.8	54.4 <sub>50.1</sub> 58.7
A	SVM	12.3 <sub>11.7</sub> 13.0	50.1 <sub>46.8</sub> 53.3	45.2 <sub>42.6</sub> 47.8	41.6 <sub>39.2</sub> 43.9
E	SVM	12.9 <sub>12.2</sub> 13.5	49.3 <sub>45.9</sub> 52.7	43.7 <sub>41.7</sub> 45.7	39.6 <sub>37.9</sub> 41.3
E	MLP	13.3 <sub>12.6</sub> 14.0	52.1 <sub>49.0</sub> 55.1	50.4 <sub>47.7</sub> 53.2	49.7 <sub>46.9</sub> 52.5

Table 9: Results of systems for Abtractivity Level classification task in *Median* dataset.

1067 Table 10 details the results obtained all *AIFs-to-*  
 1068 *AL* and *End-to-End* systems for Abstractivity Level  
 1069 regression task with *Annotators* dataset, and top-10  
 1070 configurations for *Pipeline* systems.

Type	Mthd	R-Dist ↓	RMSE ↓	MdAE ↓
A	ISVM	13.62 <sup>13.13</sup> <sub>14.11</sub>	0.88 <sup>00.85</sup> <sub>00.91</sub>	0.56 <sup>00.53</sup> <sub>00.59</sub>
A	SVM	13.77 <sup>13.30</sup> <sub>14.23</sub>	0.89 <sup>00.86</sup> <sub>00.91</sub>	0.56 <sup>00.54</sup> <sub>00.59</sub>
A	MLP	13.83 <sup>13.35</sup> <sub>14.31</sub>	0.88 <sup>00.86</sup> <sub>00.91</sub>	0.56 <sup>00.54</sup> <sub>00.61</sub>
A	LiR	13.96 <sup>13.37</sup> <sub>14.56</sub>	0.89 <sup>00.85</sup> <sub>00.93</sub>	0.60 <sup>00.56</sup> <sub>00.63</sub>
A	RnF	14.30 <sup>13.75</sup> <sub>14.85</sub>	0.90 <sup>00.87</sup> <sub>00.93</sub>	0.60 <sup>00.56</sup> <sub>00.65</sub>
P	SVM+ISVM	<b>14.73</b> <sup>14.09</sup> <sub>15.37</sub>	0.95 <sup>00.92</sup> <sub>00.98</sub>	0.59 <sup>00.56</sup> <sub>00.62</sub>
E	ISVM	14.85 <sup>14.20</sup> <sub>15.50</sub>	0.96 <sup>00.93</sup> <sub>01.00</sub>	<b>0.56</b> <sup>00.54</sup> <sub>00.61</sub>
P	ISVM+ISVM	14.88 <sup>14.27</sup> <sub>15.49</sub>	0.95 <sup>00.92</sup> <sub>00.98</sub>	0.60 <sup>00.56</sup> <sub>00.63</sub>
P	MLP+ISVM	14.97 <sup>14.43</sup> <sub>15.50</sub>	<b>0.94</b> <sup>00.91</sup> <sub>00.97</sub>	0.64 <sup>00.61</sup> <sub>00.67</sub>
P	LiR+ISVM	14.97 <sup>14.41</sup> <sub>15.52</sub>	<b>0.94</b> <sup>00.91</sup> <sub>00.97</sub>	0.63 <sup>00.60</sup> <sub>00.66</sub>
P	ISVM+SVM	14.99 <sup>14.36</sup> <sub>15.62</sub>	0.98 <sup>00.94</sup> <sub>01.01</sub>	<b>0.56</b> <sup>00.55</sup> <sub>00.60</sub>
P	LiR+SVM	15.02 <sup>14.42</sup> <sub>15.62</sub>	0.97 <sup>00.93</sup> <sub>01.00</sub>	0.59 <sup>00.56</sup> <sub>00.62</sub>
P	RnF+SVM	15.06 <sup>14.47</sup> <sub>15.65</sub>	0.97 <sup>00.94</sup> <sub>01.00</sub>	0.62 <sup>00.59</sup> <sub>00.65</sub>
P	SVM+SVM	15.08 <sup>14.45</sup> <sub>15.70</sub>	0.98 <sup>00.95</sup> <sub>01.02</sub>	0.59 <sup>00.56</sup> <sub>00.62</sub>
P	RnF+ISVM	15.09 <sup>14.53</sup> <sub>15.64</sub>	0.95 <sup>00.92</sup> <sub>00.98</sub>	0.66 <sup>00.63</sup> <sub>00.69</sub>
P	MLP+SVM	15.09 <sup>14.53</sup> <sub>15.65</sub>	0.97 <sup>00.93</sup> <sub>01.00</sub>	0.59 <sup>00.56</sup> <sub>00.62</sub>
E	SVM	15.13 <sup>14.55</sup> <sub>15.72</sub>	0.97 <sup>00.94</sup> <sub>01.01</sub>	0.61 <sup>00.59</sup> <sub>00.64</sub>
E	MLP	15.31 <sup>14.84</sup> <sub>15.78</sub>	0.94 <sup>00.92</sup> <sub>00.97</sub>	0.66 <sup>00.64</sup> <sub>00.69</sub>
E	LiR	15.40 <sup>14.76</sup> <sub>16.03</sub>	0.96 <sup>00.91</sup> <sub>01.00</sub>	0.65 <sup>00.61</sup> <sub>00.70</sub>
E	RnF	15.45 <sup>14.91</sup> <sub>16.00</sub>	0.96 <sup>00.93</sup> <sub>00.99</sub>	0.65 <sup>00.62</sup> <sub>00.68</sub>

Table 10: Results systems for Abstractivity Level regression task in *Annotators* dataset.