

A full end-to-end deep approach for detecting and classifying jaw movements from acoustic signals in grazing cattle

Mariano Ferrero^a, Leandro D. Vignolo^a, Sebastián R. Vanrell^a,
Luciano Martinez-Rau^a, José O. Chelotti^{a,b}, Julio R. Galli^c,
Leonardo L. Giovanini^a, H. Leonardo Rufiner^{a,d}

^a*Instituto de Investigación en Señales, Sistemas e Inteligencia Computacional, sinc(i),
FICH-UNL/CONICET, 3000 Santa Fe, Argentina*

^b*TERRA Teaching and Research Center, University of Liège, Gembloux Agro-Bio Tech
(ULiège-GxABT), 5030 Gembloux, Belgium*

^c*Instituto de Investigaciones en Ciencias Agrarias de Rosario, IICAR, Facultad de
Ciencias Agrarias, UNR-CONICET, Parque J.F. Villarino, S2125 Zavalla, Argentina*

^d*Laboratorio de Cibernética, Facultad de Ingeniería, Univ. Nacional de Entre Ríos, 3100
Oro Verde, Argentina*

Abstract

1 Monitoring the foraging behaviour of ruminants is a key task to improve
2 their productivity and welfare. During the last decades, several monitoring
3 approaches have been proposed based on different types of sensors such as
4 pressure-based, accelerometers and microphones. Among them, microphones
5 have been one of the most promising options because they acoustic signals
6 provide comprehensive information about the foraging behaviour. In this
7 work, a fully end-to-end deep architecture is proposed in order to perform
8 both detection and classification tasks of masticatory events in one step, re-
9 lying only on raw acoustic signals. The main benefit of this novel approach is
10 the substitution of handcrafted preprocessing and feature extraction phases
11 for a pure deep learning approach, which has shown better performance in re-
12 lated fields. Furthermore, different data augmentation techniques have been

Email address: mferrero@sinc.unl.edu.ar (Mariano Ferrero)

13 evaluated to address the data shortness for models development, typical in
14 this field. The results demonstrate that the proposed architecture achieves
15 a F1 score value of 79.82, which represents an increment close to 18% with
16 respect to other state-of-the-art algorithms. Moreover, the proposed data
17 augmentation techniques provide further performance enhancements, emerg-
18 ing as interesting alternatives in this field.

Keywords: Deep learning, data augmentation, acoustic monitoring,
precision livestock farming, ruminant foraging behaviour.

19 1. Introduction

20 Specific changes in animal behaviour are directly related to its physical
21 conditions (Frost et al., 1997), therefore tracking these changes comprises an
22 essential task of livestock management monitoring. Traditionally, it has been
23 done by manual observation, which is labour-intensive and unfeasible in some
24 practical scenarios. With the advances in communication and information
25 technologies, new automatic and non-invasive methods arose to boost data
26 collection and processing, simplifying herd management tasks (Neethirajan,
27 2020).

28 Monitoring ruminants' foraging behaviour is a critical and challenging
29 task. When long-term analyses are performed (ranging from several minutes
30 to hours), two main activities must be distinguished: rumination and graz-
31 ing. These activities are build-up on different jaw movement (JM) events:
32 bites, chews and chew-bites (Ungar et al., 2006; Milone et al., 2012). Bites
33 reflect the apprehension and severance of forage, and chews, the herbage
34 comminution. A combination of them in the same jaw movement is called

35 a chew-bite event. Monitoring the number of these events helps to provide
36 useful information regarding animal health, nutrition status, welfare and for-
37 aging activities (De Boever et al., 1990). For example, a consistent reduction
38 in rumination activity might indicate the presence of health disorders or dis-
39 eases (Calamari et al., 2014; Paudyal et al., 2018).

40 Different sources of information have been used in the last decades to
41 detect and classify JM events (Andriamandroso et al., 2016; Monteiro et al.,
42 2021). Initially, the proposed strategy was based on observation (in-situ or
43 video recordings), switches and jaw strap adjustment (Balch, 1958; Penning,
44 1983; Matsui and Okubo, 1991). This complex and fault-prone solution heav-
45 ily depends on experts and is not possible to automate it, being unfeasible
46 in large herds (Milone et al., 2009).

47 Other methods that recognise JM events rely on pressure sensors mounted
48 in a halter. The RumiWatch system (Itin and Hoch GmbH, Liestal, Switzer-
49 land) is comprised of a pressure sensor and a 3D accelerometer to gather
50 data produced during JM. This data is later analysed by a software that
51 discriminates between chews produced during rumination, chews produced
52 during feeding and grazing bites (Rombach et al., 2019). Although this sen-
53 sor reached good performance under different conditions (Ruuska et al., 2016;
54 Werner et al., 2018), their main limitation is the requirement of human inter-
55 vention for calibration, making infeasible its use in commercial farms (Riaboff
56 et al., 2022). Additionally, several practical issues have been reported in the
57 use of halters (Nydegger et al., 2011) such as frequent damage when applied
58 in loose housing systems.

59 On the other hand, diverse motion sensors located in different places of the

60 animal's body have been used to determine long-term activities (rumination,
61 grazing, resting, among others) rather than JM events (Fogarty et al., 2020;
62 Balasso et al., 2021; Riaboff et al., 2022).

63 Bite events count has been addressed using pattern matching techniques
64 from 1D accelerometer (Tani et al., 2013), 3D accelerometer (Oudshoorn
65 et al., 2013; Giovanetti et al., 2017) and inertial measurement unit (Andria-
66 mandroso et al., 2015). Despite the fact that motion sensors provide inter-
67 esting options to automatically count feeding JM (low sampling frequency
68 and comprehensive data), the distinction between different types of events
69 represents a challenging task from these signals and proper validation on di-
70 verse pasture and larger duration trials is still required (Ding et al., 2022).
71 The sensitivity of this kind of sensors might introduce errors and misclassi-
72 fications due to unrelated movements with JM events (ear wiggling or head
73 turns). Furthermore, position displacements of the motion sensor affect the
74 JM event recognition, and they are difficult to prevent in free-ranging condi-
75 tions (Kammaing et al., 2018; Li et al., 2021a).

76 Acoustic sensors are useful for the recognition of JM events in free-ranging
77 environments. The use of microphones allows for capturing the sounds pro-
78 duced by the teeth and propagated through the bones, cavities and soft tis-
79 sues of the cattle's head. The analysis of these signals is a difficult task due
80 to the presence of environmental sounds (noises) and the high computational
81 requirements. Beyond that, they are usually preferred over pressure and
82 movement sensors because the acoustic signals capture more information in
83 order to perform JM events classification (Ungar et al., 2006; Martinez-Rau
84 et al., 2022). Milone et al. (2012) developed a computational demanding

85 method to detect and classify JM events using hidden Markov models on
86 spectral-domain features. Navon et al. (2013) proposed a machine learning
87 approach to separate true events (without specific classification) from back-
88 ground noise and silence. Chelotti et al. (2016) proposed the Chew-Bite
89 Real-Time Algorithm, which defined a sequential system for detecting and
90 classifying chews, bites and chew-bites using heuristic rules and temporal fea-
91 tures. In a later work, searching for better results, the same authors proposed
92 a system based on machine learning called Chew-Bite Intelligent Algorithm
93 (CBIA) (Chelotti et al., 2018). Recently, Martinez-Rau et al. (2022), pro-
94 posed an algorithm for robust recognition of JM events called Chew-Bite
95 Energy Based Algorithm. It is capable of discriminating four event types:
96 bites, chew-bites, rumination chews and grazing chews.

97 Automatic detection and classification systems based on sound analysis
98 usually perform a preprocessing stage (e.g., to improve signal-to-noise-ratio)
99 and then execute some sort of feature extraction to feed data into the classifi-
100 cation models. The lack of an end-to-end solution introduces several potential
101 troubles, such as dependency on specific sound recording systems and config-
102 uration, as well as difficulties to exploit potentially valuable information not
103 encoded in manually created features. Li et al. (2021c) introduced a compar-
104 ison of several deep learning (DL) architectures to classify JM events using
105 a preprocessing phase where frequency-domain representations are extracted
106 from raw signals. The complete workflow proposed by these authors, to gen-
107 erate the inputs of neural networks models includes the following steps: back-
108 ground noise removal using a band-stop filter, uninformative data removal
109 based on manually created thresholds and Mel-frequency cepstral coefficients

110 calculation. Compared with traditional machine learning techniques, the use
111 of DL models brings the opportunity to automatically discover patterns and
112 features from data at the expense of higher computational costs.

113 Based on the analysis of previous research it is possible to state that DL
114 models have been used only to classify JM events. Therefore, the application
115 of DL models to perform JM events recognition (which involves JM events
116 detection and the posterior classification of them), has not been explored
117 yet. Additionally, the rest of the traditional alternatives (such as the CBIA
118 system) heavily depend on manual feature extraction methods and arbitrarily
119 defined pre-processing steps. Promising results presented by Li et al. (2021c)
120 highly motivate the study of DL architectures to tackle the limitation of JM
121 events recognition.

122 In this paper, a truly end-to-end approach is proposed to process raw
123 audio signals toward the detection and the classification of JM events (bite,
124 chew and chew-bite). The proposed DL strategy combines the power of con-
125 volutional networks for feature learning with the time modeling capabilities
126 of recurrent units, to implement **detection and classification tasks in one**
127 **step**. Several architectures have been explored and compared to point out
128 the benefits and limitations of the proposed approach. Additionally, different
129 data augmentation techniques have been evaluated to improve the generali-
130 sation capabilities of the proposed approach. Experimental results show the
131 benefits of the application of the proposed deep architectures over traditional
132 machine learning approaches. The main contributions of this paper are the
133 following: a) a novel deep-learning model that combines convolutional and
134 recurrent neural networks is presented. It automatically learns the features

135 representations and the temporal dependencies between JM events from raw
136 audio signals. b) The proposed model is able of solving the JM events detec-
137 tion and classification tasks in one step from raw from acoustic signals; and
138 finally c) different data augmentation techniques were analysed to undertake
139 the data-shortness problem.

140 2. Material and methods

141 In this article, a novel deep-learning architecture called **Deep sound** is
142 proposed. It is based on the combination of two types of neural networks:
143 Convolutional Neural Networks (CNN) (Lecun et al., 1998) and Recurrent
144 Neural Networks (RNN) (Rumelhart et al., 1986). In the following sections,
145 a brief introduction to these architectures is provided. Then, a detailed
146 description of the proposed method is presented.

147 2.1. CNN and RNN

148 Convolutional Neural Networks (CNN) (Lecun et al., 1998) are one of
149 the most widely used architectures for classification problems where input
150 data comes from unstructured sources - images (Kokalis et al., 2020) or au-
151 dio (Ramirez et al., 2022), for example. They are usually composed by
152 several convolutions layers, each one containing one or more filters. In the
153 learning stage, filters' weights (used in traditional convolution mathematical
154 operations) are adapted in order to approximate outputs using optimisation
155 strategies like stochastic gradient descent or back-propagation (Rumelhart
156 et al., 1986). By doing this, the layers are capable of learning different high
157 and low-level patterns without domain knowledge supplied.

158 In CNN, convolutional layers are used in combination with pooling, batch
159 normalisation and dense layers. Pooling layers apply simple mathematical
160 operations (such as maximum extraction) in order to reduce dimensionality,
161 and they are commonly used after convolutional layers. On the other hand,
162 batch normalisation layers scale the inputs, to the desired values, to accel-
163 erate the training process. Finally, dense layers correspond to a flat set of
164 hidden neurons fully connected (FNN) with the outputs of previous layers,
165 providing to the CNN with the ability to adapt the effect of intermediate
166 representations, learned by convolutions, on the output. The relation be-
167 tween convolution with other layers is created using a flattening operation,
168 which transforms the output of convolution layers into a vector. An impor-
169 tant operation used in these layers (except for batch normalisation) is called
170 **drop-out**, which introduces random crops between layer connections during
171 the training phase to avoid model over-fitting (Hinton et al., 2012).

172 Recurrent Neural Networks (RNN) (Rumelhart et al., 1986) are broadly
173 used in a wide variety of problems involving temporal sequences (Lim et al.,
174 2019; Li et al., 2021b). RNN connects layer outputs as inputs to the same
175 layer, enabling temporal data flow more efficiently across the network. More
176 sophisticated architectures have been developed in recent years to overcome
177 some RNN limitations. Gated Recurrent Units (GRU) are composed of sev-
178 eral neurons called **cells**, each one uses two different gates: reset and update
179 (Cho et al., 2014). These gates, tuned during the training process, allow
180 every neuron to control the trade-off between how much information is used
181 from previous and current states. GRU networks are composed of several
182 GRU cells placed sequentially. A variation of a RNN proposed by Schuster

183 and Paliwal (1997) is called Bidirectional RNN. This network introduces two
184 identically RNN in terms of architecture, one trained with time sequences for-
185 wards and the other one with the same sequences backward, both connected
186 to the next layer of the network. Specifically, bidirectional GRU (BGRU)
187 achieved very promising results in sound events detection (Lu et al., 2018;
188 Meng et al., 2022) and classification (Zhu et al., 2020).

189 *2.2. Deep sound*

190 Different variations of several deep architectures were studied for this
191 problem, based on previous research in related fields (Khamees et al., 2021;
192 Bahmei et al., 2022; Petmezas et al., 2022). The alternatives were evalu-
193 ated from a theoretical perspective and the most promising ones were im-
194 plemented. Thus, a hybrid one-dimensional (1D) CNN-BGRU network ar-
195 chitecture is proposed, named **Deep sound**. To the best of the authors'
196 knowledge, this represents the first deep end-to-end approximation to the
197 problem of JM events detection and recognition from acoustic signals. The
198 network receives the sound windows extracted from the original audio files
199 without any prior preprocessing or feature extraction phase, and classifies
200 them into one of four possible classes: chew, bite, chew-bite and no-event.
201 Therefore, the proposed method tackles the problems of JM event detection
202 and classification at the same time.

203 The proposed model structure is given by: an input layer and several hid-
204 den layers distributed in three main blocks corresponding to CNN, BGRU,
205 and FNN. An overall schematic of the proposed model is presented in Fig-
206 ure 1(a), while a detailed description of the architecture is showed in Fig-
207 ure 1(b). The first part of Figure 1(b) represents the CNN block of the model,

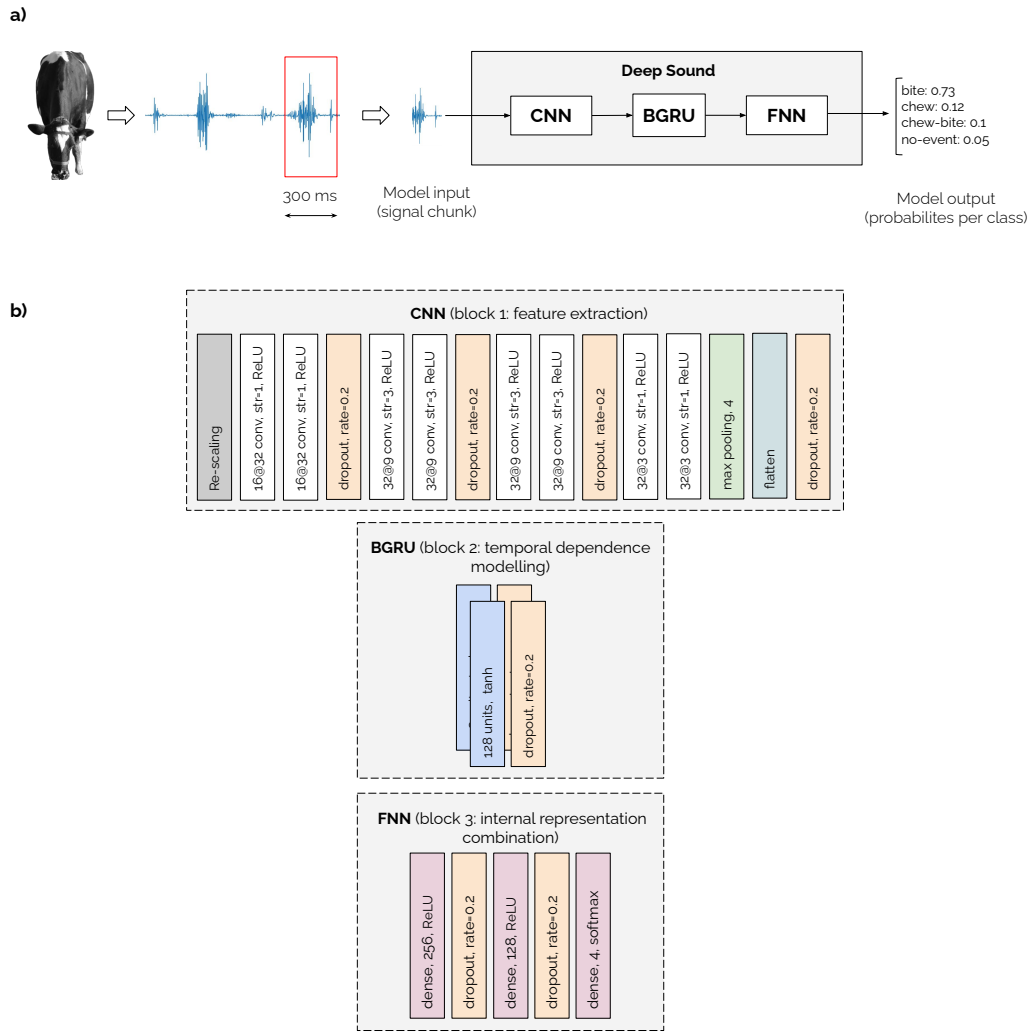


Figure 1: The overall proposed method architecture. a) Input signals correspond to audio chunks extracted using fixed-length time windows and passed through the CNN (first block) to automatically extract features. The output of this block is passed to the bidirectional GRU to capture temporal dependencies in data. Finally, the output of the second block is fed into the FNN block, combining information in dense layers, and predicts class probabilities for each input sample. b) Specification of layers in each block, including the number of filters or units, filter size (for convolutional layers), and activation functions.

208 which is a combination of 1D convolutional layers, dropout operations, and
209 max pooling layers. This way, the network is capable of extracting low- and
210 high-level features from audio chunks and performing dimensionality reduc-
211 tion at the same time. At the beginning of this block, a re-scaling layer
212 adapts the range of input values for implementation purposes. A flatten op-
213 eration is also used to create a raw vector from the last convolutional layer. A
214 complete definition of layer configurations, such as number of filters and filter
215 sizes, is provided in the figure. The second block in Figure 1(b) introduces a
216 recurrent network, composed of a BGRU layer with 128 cells. The purpose
217 of this block is to capture time dependencies in the data. The last block
218 of the network implements a typical FNN with three dense layers and two
219 dropout operations. Blocks one and three are placed into time-distributed
220 wrappers, allowing the same layers to be applied to each window of the in-
221 put signals. This means that the same set of connection weights is trained
222 and used in these blocks for every time window. All convolutional layers
223 use the activation function rectified linear unit (ReLU), whilst the cells of
224 the BGRU use hyperbolic tangent and sigmoid. The first and second dense
225 layers perform both ReLU, and the last dense layer uses the soft-max func-
226 tion for classification. All layers (convolutional, recurrent and dense) use the
227 Xavier initialisation method (Glorot and Bengio, 2010) and bias terms were
228 initialised to zero.

229 The main limitations of the proposed method are: a) a considerable
230 amount of labelled data is needed for training, b) the interpretability of the
231 method and its outputs is limited (Arrieta et al., 2020; Hoxhallari, 2022), and
232 c) a considerable amount of processing is required in the inference phase.

233 *2.3. Acoustic dataset*

234 *2.3.1. Original dataset*

235 The data used in this work is one of the first open datasets in this field of
236 study (Vanrell et al., 2020). The fieldwork to obtain this dataset took place
237 at the Campo Experimental J.F. Villarino, Facultad de Ciencias Agrarias,
238 Universidad Nacional de Rosario, Zavalla, Argentina. The recordings include
239 sounds produced by dairy cows in individual grazing sessions conducted over
240 a 5-day period. Microphones used to record audio signals (Nady 151 VR,
241 Nady Systems, Oakland, CA, USA) were located on the cow’s forehead and
242 covered with rubber foam. Further details about experimental design could
243 be found in the dataset article (Vanrell et al., 2020).

244 A total of 52 raw audio signals (WAV audio files, mono, 16-bits, 22.05
245 kHz) are available ¹. A summary of the dataset contents is presented in Ta-
246 ble 1. Each audio signal consists of sequences of JM events – bites, chews,
247 and chew-bites – separated by silence (ranging from 19 to 152 s, average du-
248 ration 62.76 ± 28.61 s). Two different experts in ruminant foraging behaviour
249 independently performed the identification of each JM (including event la-
250 bel, start, and end time) by analysing videotapes and sounds at the same
251 time. Agreement results were 100% for bites, 98.2% for chews, and 99.1%
252 for chew-bites. There were 2.7% of insertions and 0.9% of deletions. Thus,
253 the total segmentation and classification accuracy was 93.6%. Both experts
254 worked together to achieve a final decision in case of disagreement.

¹Direct URL to data: <https://github.com/sinc-lab/dataset-jaw-movements>

Table 1: Summary of audio files grouped by pasture and height.

Pasture	Height	Chews	Bites	Chew-Bites	Overall duration
Alfalfa	Tall	416	148	322	14 min 26 s
Alfalfa	Short	260	179	123	12 min 42 s
Fescue	Tall	487	100	238	14 min 03 s
Fescue	Short	454	94	217	13 min 13 s
Total		1617 (53%)	521 (17%)	900 (30%)	54 min 24 s

255 *2.3.2. Data preparation*

256 Since the delimitation of most of the labels in the original dataset was
 257 inaccurate with respect to the actual JM events, an improvement to label
 258 bounds has been proposed in the present work. Conducting a visual inspec-
 259 tion of original signals and labels, it is possible to notice that there is not a
 260 perfect time delimitation between JM events presence and timestamps. Fig-
 261 ure 2 shows some examples where over estimations of JM events duration
 262 are introduced. To tackle this situation, time event delimiters have been
 263 adapted using a label erosion method based on signal envelope computation
 264 and selected thresholds. The events start timestamp was moved to the po-
 265 sition where the signal envelope reaches a certain threshold; similarly, this
 266 process was repeated in the opposite direction with the event end timestamp,
 267 generating a time shift respecting the original label.

268 The threshold is defined as follows: after JM event envelope calculation,
 269 the maximum value is obtained and multiplied by a factor adapted to the
 270 differences between event characteristics. Table 2 introduces start and end
 271 factors applied to different event classes.

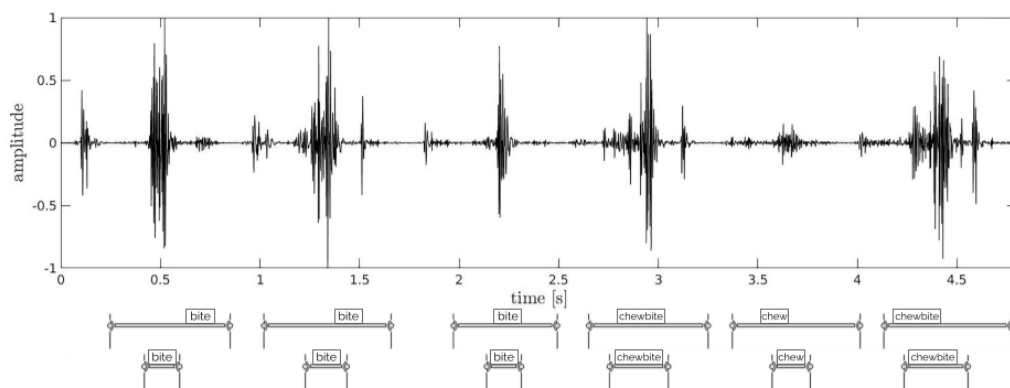


Figure 2: Visual comparison of an example of a signal with original (top) and eroded (bottom) labels with time delimiters (timescale on the top is expressed in seconds).

Table 2: Scale factors applied to maximum values extracted from the signal envelope to define threshold calculation.

JM event type	Start factor	End factor
Bite	0.4	0.4
Chew	0.5	0.5
Chew-Bite	0.15	0.4

272 Original audio signals have been recorded at 22.05 kHz. In order to reduce
273 dimensionality and computational costs, all files were downsampled to 6 kHz.
274 In addition to this, original audio signals were divided into small chunks
275 of data using sequentially ordered windows. Different window sizes have
276 been evaluated during the initial experimentation, considering the average
277 duration of JM events, and the value of 300ms produced the best results,
278 with a hop length of 150 ms. The average duration of the JM events is 330
279 ms (± 150 ms), which means that two consequent windows might be needed
280 to represent one JM event. To assign a label to a particular signal window, a
281 minimum overlapping of 40% with a JM event label is required, guaranteeing
282 that if only a small part of a window corresponds to a JM event of interest
283 (bite, chew or chew-bite) it is tagged as 'no-event'.

284 2.3.3. Data augmentation

285 A distinctive characteristic of the proposed approach is the number of pa-
286 rameters to be learned or tuned during the training process. Consequently,
287 the use of a small amount of data may lead to overfitting. In the context of
288 precision livestock farming, and JM events recognition in particular, getting
289 more annotated signals requires great effort and resources. To overcome this
290 problem, data augmentation techniques are traditionally employed to artifi-
291 cially create synthetic samples from original ones (Nanni et al., 2021; Bahmei
292 et al., 2022). Despite that data augmentation is well-known for image-related
293 problems (Shorten and Khoshgoftaar, 2019), custom techniques are usually
294 required when working with audio signals.

295 When new samples are created from existing data, two facts should be
296 considered: *i*) the types of perturbations applied on original data to create a

297 different one, but still usable synthetic audio signal (named here **augmen-**
298 **tation technique**), and *ii*) how to apply them to every training sample
299 (**augmentation protocol**). Several augmentation techniques have been ex-
300 plored in early experimentation (including but not limited to loop, pitch shift,
301 time stretch and percussive). Finally, six data augmentation techniques were
302 selected:

- 303 • Resynthesis by Linear Predictor Coefficients (LPC): given an input
304 signal, the LPC is estimated, randomly perturbed, and finally used to
305 generate a new signal using a resynthesis process.
- 306 • Reverse: a copy is created from original values by doing a backward
307 pass.
- 308 • Random crop: randomly pick a very small fraction (1%) of continuous
309 values from the input signal and turn them to zero.
- 310 • Background noise: add white noise to the original signal, using a signal-
311 to-noise ratio of 10 dB.
- 312 • Amplitude change: increase or decrease signal amplitude by a certain
313 decibel amount. Positive values stand for increases, while negative
314 stands for amplitude decrease.
- 315 • Frequency filters: apply a second-order Butterworth high-pass or low-
316 pass filter to the input signal. The high-pass and low-pass filters have
317 a cut-off frequency of 500 Hz and 100 Hz, respectively.

318 On the other hand, two different augmentation protocols were tested:

- 319 • Random: pick one augmentation technique and use it to generate a
320 synthetic signal.
- 321 • Serial: create a pipeline serialising all defined augmentation techniques
322 in order to apply them one by one. This way the input to the first tech-
323 nique is the original audio signal and its output is fed to the subsequent
324 technique.

325 During experimentation, three synthetic signals were created from every
326 single input sample when defining an augmentation protocol. These values
327 were selected in order to explore the effect of this component without signif-
328 icantly affecting the computational cost.

329 2.4. Experimentation methodology

330 2.4.1. Model selection approach

331 For all experiments, the models were evaluated using 10-fold cross-validation
332 (CV). Every fold contains 5 or 6 input files, randomly selected from the total
333 of 52 available. In this way, every input file was included in only one fold.
334 In addition to this, 20% of the 9 folds used for training on every iteration
335 were reserved for validation. The assignment of sound files to the train and
336 test sets in each fold was fixed across different experiments. The number of
337 windows in test sets was 2168 ± 360 (proportion per class: $5 \pm 1\%$ bites -
338 $18 \pm 1\%$ chews - $14 \pm 4\%$ chew-bites - $63 \pm 4\%$ no-event). The number of
339 windows in train and validation sets changed from one experiment to another
340 due to the use of different data augmentation configurations. The training
341 samples were weighted in order to tackle classes imbalance according to the
342 following expression:

$$cw_{ic} = n_{max}/n_c, \quad (1)$$

343 where cw_{ic} is the class weight of instance i of class c , n_{max} is the number
 344 of instances of the majority class and n_c is the number of instances of class
 345 c . Finally, the experiments were set-up with a total of 1500 epochs with
 346 early stopping (50 epochs tolerance), Adam (Kingma and Ba, 2014) as the
 347 optimizer, the batch size was fixed to 10, 0.001 as the learning rate, and
 348 categorical cross entropy as loss function. Default values were used for the
 349 remaining parameters.

350 2.4.2. Evaluation metrics

351 The dynamical problem of simultaneous detection and classification of JM
 352 events using raw audio signals is substantially different from the approach of
 353 dividing the problem into JM event detection and subsequent classification
 354 based on previously detected events (Chelotti et al., 2018; Martinez-Rau
 355 et al., 2022). In the former, the temporal component plays a very important
 356 role, since the need to properly detect JM event’s onsets and offsets affects
 357 the results of the classification. Based on this, the generation of a model
 358 that deals with detecting and classifying events at once requires the use
 359 of a validation mechanism that is capable of considering aspects related to
 360 temporality, as well as predicted labels accuracy.

361 To evaluate JM events detection and classification performances, the
 362 `sed_eval` standardised toolbox was used (Mesaros et al., 2021). It is a trans-
 363 parent and broad library to evaluate sound event recogniser systems. The
 364 toolbox was designed for the task of sound event recognition, which involves

365 locating and classifying sounds in audio recordings, estimating onset and off-
 366 set for distinct sound event instances and providing a textual descriptor for
 367 each. This matches the task presented in this work, where sound JM events
 368 classes are chew, bite and chew-bite. A temporal tolerance (collar) of 300 ms
 369 was used. This value was determined based on preliminary experimentation
 370 considering two main aspects: 1) the collar should be smaller than the aver-
 371 age event duration (330 ms) in order to ensure overlap between reference and
 372 predicted window. 2) it should avoid undesired overlap between two adjacent
 373 events (with an average separation between two adjacent events of 726 ms).
 374 The selected value meets both criteria.

375 With the use of `sed_eval` toolbox, a reference JM event is correctly de-
 376 tected if two conditions are met: *i)* The start timestamp of the predicted JM
 377 event is located in the interval defined by reference onset \pm tolerance value.
 378 *ii)* The end timestamp of the predicted JM event is located in the interval
 379 defined by reference offset \pm tolerance value. Figure 3 introduces a graphical
 380 representation of how this toolbox works.

381 Based on before mentioned evaluation toolbox, several well-known metrics
 382 have been used:

$$precision = \frac{TP}{TP + FP},$$

$$recall = \frac{TP}{TP + FN},$$

$$F1\ score = \frac{2 * precision * recall}{precision + recall},$$

$$error\ rate = \frac{S + D + I}{N}$$

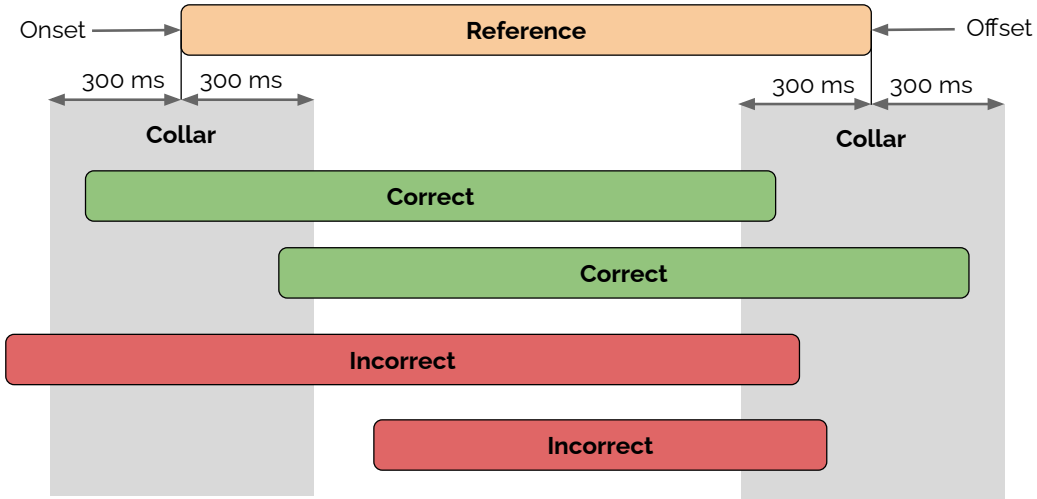


Figure 3: Illustration based on Mesaros et al. (2021) where two correct and two incorrect predicted JM events are presented, compared with a reference JM event using a tolerance value of 300 ms.

386 where TP denotes true positive, FP false positive, FN false negative, S
 387 substitutions (correct detected JM events in system output but incorrectly
 388 labelled), I insertions (detected events from system output which do not
 389 exist in the ground truth), D deletions (ground truth events which are not
 390 detected) and N is the total number of reference events. Due to the presence
 391 of class imbalance in the original dataset, JM events distributions are taken
 392 into account to calculate average final results. When using this approach for
 393 metrics calculation micro averages were computed (Sokolova and Lapalme,
 394 2009), which means that TP , FP and FN are calculated by summing up
 395 samples through all classes. For example, the term TP is finally expressed
 396 by $TP_c + TP_{cb} + TP_b$, representing the amount of TP for chews, chew-bites
 397 and bites, respectively.

398 *2.5. Experimental setup*

399 The design and implementation of the proposed model were developed
400 using Python 3.6.2 and TensorFlow-GPU 2.6.2. Different utilities from the
401 Python library scikit-learn 0.24.2 have been used, such as label encoders
402 and k-fold extraction. Augly (Papakipos and Bitton, 2022), a Python data
403 augmentation library, was used to apply some of the previously mentioned
404 augmentation techniques (background noise, amplitude change and frequency
405 filters). Experiments were performed using an Intel[®] Core[™] i7-8700 3.20GHz
406 CPU, 64 GB RAM and 24 GB NVIDIA GeForce RTX 3090 GPU. A Titan
407 XP GPU was also used for model exploration, preliminary experimentation
408 and hyperparameter tuning.

409 **3. Results**

410 During the optimisation process, a total of 39 experiments were tested,
411 aiming to find the best model architecture configuration considering varia-
412 tions in the CNN part of the model (block 1 in Figure 1). The most promising
413 and standard hyper-parameters combinations (such as the number of layers,
414 number of filters, and dimension of filters) have been considered for this
415 exploration. All experiments used the 10-fold CV method described in Sec-
416 tion 2.4.1. Layers configuration from most representative experiments are
417 described in Figure 4, and their respective recognition results are presented
418 in Table 3. In terms of performance, architecture (c) exhibited the high-
419 est F1 score value. Moreover, this model also reached the lowest error rate.
420 Therefore, it is possible to establish that architecture (c) configures the best
421 combination explored, considering numbers of layers, number of filters and

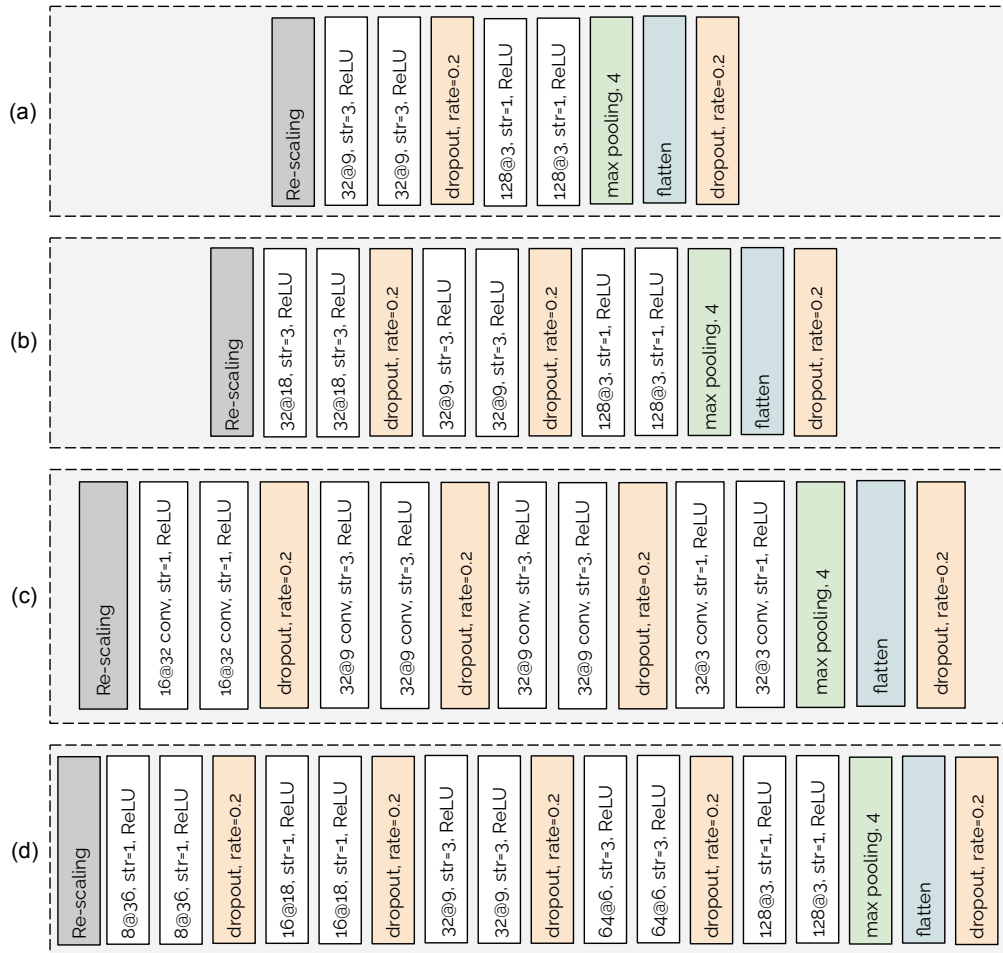


Figure 4: Different CNN architectures used for exploration. Convolution layers definition consist of number of filters, filter size and stride. No padding method was used.

422 filter dimensions.

423 As described previously, the proposed model is composed of three blocks
424 with different types of layers. Table 4 exhibits the performance of the pro-
425 posed model without using the RNN (block 2 in Figure 1). It can be seen
426 that providing the capacity to capture temporal relationships in acoustic
427 sequences gives a significant advantage to the network.

428 In addition to the optimisation of model hyperparameters, an exploration

Table 3: Recognition results of the proposed model for different layers architectures on the CNN block. For every experiment, average values and standard deviation of 10-folds CV are presented.

	Precision \uparrow	Recall \uparrow	F1 score \uparrow	Error rate \downarrow	Deletion \downarrow	Insertion \downarrow
(a)	63.13 \pm 6.53	79.81 \pm 6.06	70.45 \pm 6.26	0.54 \pm 0.12	0.07 \pm 0.03	0.34 \pm 0.07
(b)	71.91 \pm 5.26	85.77 \pm 3.37	78.19 \pm 4.33	0.39 \pm 0.08	0.05 \pm 0.02	0.25 \pm 0.06
(c)	73.72 \pm 4.92	87.16 \pm 2.74	79.82 \pm 3.70	0.37 \pm 0.08	0.05 \pm 0.01	0.24 \pm 0.07
(d)	73.38 \pm 5.30	85.92 \pm 3.81	79.12 \pm 4.46	0.37 \pm 0.09	0.06 \pm 0.02	0.23 \pm 0.06

Table 4: Evaluation of the impact of the RNN block in the proposed model. For each experiment, the average and the standard deviation of 10-fold CV are presented.

	Precision \uparrow	Recall \uparrow	F1 score \uparrow	Error rate \downarrow	Deletion \downarrow	Insertion \downarrow
Deep sound	73.72 \pm 4.92	87.16 \pm 2.74	79.82 \pm 3.70	0.37 \pm 0.08	0.05 \pm 0.01	0.24 \pm 0.07
Deep sound (no RNN)	48.77 \pm 3.89	82.55 \pm 3.64	61.26 \pm 3.79	0.95 \pm 0.14	0.07 \pm 0.03	0.77 \pm 0.12

429 of the impact of using several data augmentation techniques and protocols
 430 were carried out using the proposed Deep sound (c) architecture. Table 5
 431 introduces the results of different experiments using isolated augmentation
 432 techniques (in order to measure the individual impact) and combining many
 433 of them at the same time with a particular augmentation protocol. The
 434 protocol combined the three best individual techniques based on its F1 score
 435 (background noise, random crop and amplitude (+2 dB)) to form a top 3
 436 augmentation technique. This combination has been tested using serial and
 437 random protocols. The highest F1 score (p=0.006; Wilcoxon signed-rank
 438 test) (Wilcoxon, 1945) was reported using the top 3 augmentation techniques
 439 with serial augmentation protocol.

440 Finally, a contrast between the proposed model and other state-of-the-art

Table 5: Results of the proposed model using different augmentation techniques and protocols. For each experiment, the average and the standard deviation of 10-fold CV are presented. The number of copies generated per original sample was fixed to three.

Augmentation technique	Augmentation protocol	Precision \uparrow	Recall \uparrow	F1 score \uparrow	Error rate \downarrow
No augmentation	-	73.72 \pm 4.92	87.16 \pm 2.74	79.82 \pm 3.69	0.37 \pm 0.08
LPC	-	71.88 \pm 4.72	86.67 \pm 2.64	78.54 \pm 3.69	0.40 \pm 0.08
Background noise	-	76.83 \pm 5.61	85.71 \pm 3.46	80.96 \pm 4.37	0.32 \pm 0.09
Random crop	-	77.28 \pm 7.72	86.31 \pm 3.72	81.43 \pm 5.63	0.32 \pm 0.12
Amplitude (+2 dB)	-	76.14 \pm 5.33	86.60 \pm 3.89	80.98 \pm 4.37	0.33 \pm 0.08
Amplitude (-2 dB)	-	74.24 \pm 6.45	86.18 \pm 3.33	79.68 \pm 4.78	0.37 \pm 0.10
High-pass filter	-	70.63 \pm 5.57	85.25 \pm 3.82	77.19 \pm 4.59	0.42 \pm 0.09
Low-pass filter	-	66.64 \pm 8.37	83.80 \pm 4.99	74.09 \pm 6.83	0.50 \pm 0.17
Reverse	-	72.90 \pm 5.91	86.78 \pm 2.61	79.16 \pm 4.38	0.39 \pm 0.09
Top 3	Serial	78.39 \pm 4.09	86.60 \pm 3.08	82.27 \pm 3.42	0.29 \pm 0.06
Top 3	Random	77.04 \pm 5.45	87.06 \pm 3.19	81.67 \pm 3.99	0.32 \pm 0.08

441 methods has been carried out. In particular, the algorithm called Chew-Bite
442 Intelligent Algorithm (CBIA) (Chelotti et al., 2018) and an implementation
443 of the ResNet proposed by Hershey et al. (2017) for raw audio classifica-
444 tion were compared using the same evaluation toolbox and metrics. The
445 CBIA method was selected because it offers the best results of state-of-the-
446 art in the detection and classification of JM events problem (unlike Li et al.
447 (2021c), where only classification is performed) for chew, bite and chew-bite
448 labels. Moreover, as the authors mention in their work, the Li et al. (2021c)
449 proposal does not offer improvements in terms of classification rates with
450 respect to Chelotti et al. (2018) approach. The ResNet architecture was
451 selected because it is one of the best well-known DL models proposed for
452 image classification and reached the best results for audio classification tasks
453 (Hershey et al., 2017) among other DL models (such as VGG (Simonyan and

Table 6: Comparison between the proposed method and other state-of-the-art algorithms, CBIA and ResNet architecture.

	Precision \uparrow	Recall \uparrow	F1 score \uparrow	Error rate \downarrow	Deletion \downarrow	Insertion \downarrow
Deep sound	78.39 \pm 4.09	86.60 \pm 3.08	82.27 \pm 3.42	0.29 \pm 0.06	0.06 \pm 0.02	0.17 \pm 0.05
CBIA	68.69 \pm 7.56	70.30 \pm 7.92	69.43 \pm 7.52	0.42 \pm 0.11	0.10 \pm 0.05	0.12 \pm 0.06
ResNet audio	43.99 \pm 12.96	54.99 \pm 23.35	47.9 \pm 17.16	0.97 \pm 0.27	0.3 \pm 0.21	0.52 \pm 0.2

454 Zisserman, 2014), Inception (Szegedy et al., 2016) or AlexNet (Krizhevsky
 455 et al., 2017)).

456 The results of this comparison are presented in Table 6 and separated by
 457 class in Table 7. Deep sound refers to the best architecture configuration
 458 (architecture (c)), trained using top 3 (background noise, random crop and
 459 amplitude increase +2 dB) serial augmentation protocol. It can be seen that
 460 there is a significant improvement using the proposed algorithm ($p=0.002$
 461 based on F1 score; Wilcoxon signed-rank test) (Wilcoxon, 1945). Despite
 462 this, results from all methods are higher for chew events, probably related to
 463 the fact that this is the most predominant class. Regarding deletion metric,
 464 the proposed algorithm increases the number of ground truth events detected.
 465 However, CBIA presents a smaller number of insertions than the proposed
 466 algorithm.

467 Finally, a summary of the different approaches is introduced in Figure 5.
 468 In terms of F1 score and precision, the proposed architecture (Deep sound)
 469 using augmentation techniques obtained the best results, whereas ResNet
 470 architecture led to the lowest value. On the other hand, based on the re-
 471 call metric, the proposed architecture without augmentation techniques pre-
 472 sented the best results and ResNet produced the worst. It is possible to note

Table 7: Class based results obtained for the proposed architecture and other state-of-the-art algorithms, CBIA and ResNet architecture.

	Class	Precision \uparrow	Recall \uparrow	F1 score \uparrow
Deep sound	Bite	73.59 ± 8.49	76.10 ± 9.16	74.27 ± 6.52
	Chew	82.56 ± 6.32	90.61 ± 3.58	86.33 ± 4.78
	Chew-Bite	73.81 ± 8.40	86.53 ± 4.38	79.31 ± 5.24
CBIA	Bite	48.77 ± 10.72	66.41 ± 10.37	55.06 ± 7.48
	Chew	77.30 ± 6.59	76.69 ± 5.72	76.77 ± 4.60
	Chew-Bite	70.77 ± 15.06	60.78 ± 18.09	63.74 ± 16.65
ResNet audio	Bite	36.72 ± 20.8	55.18 ± 23.95	42.6 ± 20.7
	Chew	51.31 ± 26.02	52.6 ± 34.18	48.91 ± 28.54
	Chew-Bite	41.62 ± 12.97	62.94 ± 20.09	46.87 ± 11.95

473 that ResNet also exhibited higher deviations in all presented metrics.

474 4. Discussion

475 4.1. End-to-end model architecture

476 Based on the presented results, the use of a deep end-to-end approach
477 provides the model the capacity to learn relevant internal representations
478 starting from raw signals. Manual feature computation and extraction are
479 difficult tasks, which involve a deep understanding of the studied phenomena
480 as well as the capacity to apply that knowledge properly. This limitation
481 is overcome in the proposed model, resulting in a significant improvement
482 compared with traditional machine learning algorithms. It is important to
483 highlight that the use of recurrent layers introduces a substantial benefit to
484 the model architecture. The use of different gates allows these layers to learn

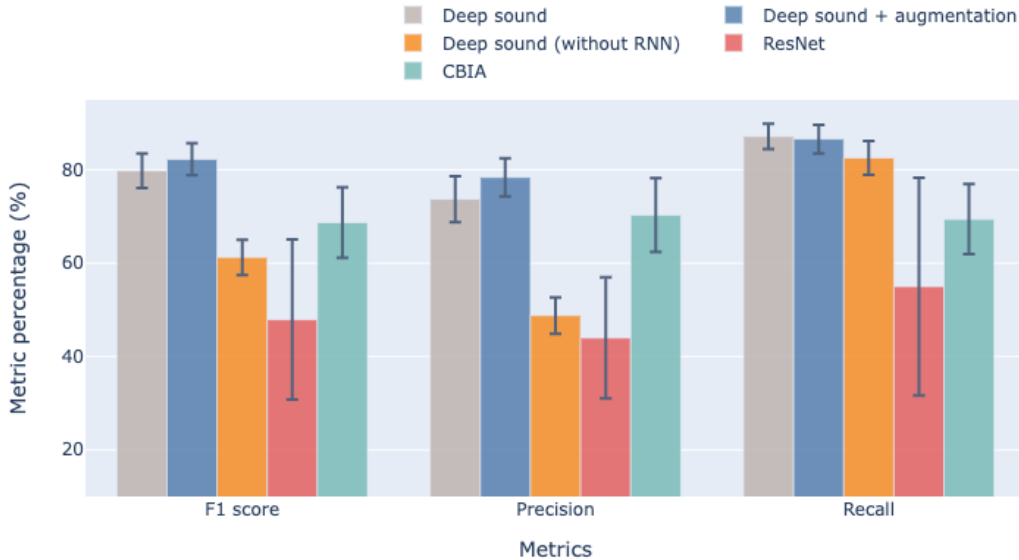


Figure 5: Overall comparison of the results obtained by the most relevant experiments of each of the presented approaches.

485 how much information to incorporate into their internal memory regarding
 486 new events and how much to remember from previous events. A positive
 487 impact seems reasonable based on this, attending to ruminant foraging be-
 488 haviour activities, in which sometimes a single bite is followed by a sequence
 489 of chew and chew-bite events during grazing.

490 Regarding the model architecture, the results suggest that the use of sev-
 491 eral layers is advantageous. When using a reduced number of convolutional
 492 layers (less than 6), the recognition performance of the network is remarkably
 493 damaged. In contrast, when using at least 6 convolutional layers the model
 494 performance seems to approach similar levels. A possible explanation of this
 495 fact is that the model requires a minimum number of layers in order to extract
 496 a relevant representation from data. In terms of the number of parameters,
 497 the model architecture presented in Figure 4 (c) uses 320,229. This value

498 probably represents a considerable increment compared to other traditional
499 methods. However, the use of convolutional layers prevents a bigger increase
500 in this number with respect to other neural network architectures, which use
501 mainly dense layers.

502 The evaluation with different data folds shows a considerable level of de-
503 viation in the performance metrics. This effect might be due to the fact that
504 several signals are particularly different from the rest in terms of duration
505 (shorter) and JM events distribution (most of the present events correspond
506 to the same class along the signal). The recognition performance decreased
507 on those signals in all performed experiments.

508 *4.2. Effect of learning from synthetic data*

509 In order to increase the size of the dataset available for training in each
510 fold, eight different data augmentation techniques were proposed and anal-
511 ysed (Table 5). Results showed that a subset of them allowed the model to
512 improve the recognition performance in terms of F1 score. When analysing
513 precision and recall separately, it is possible to note that introducing syn-
514 thetic data to the training process reduces the number of detected events
515 in general. Despite this, for some techniques there was an improvement in
516 the precision of predictions. The results highlight the importance of using
517 augmentation techniques to increase the generalisation capacity of the model.

518 Some individual techniques showed a positive impact on the performance,
519 while others showed no impact or even a negative impact. The techniques of
520 both low- and high-pass filters and reverse degraded the performance com-
521 pared to the no augmentation approach. In contrast, when adding back-
522 ground noise or random crops, the model presented improvements regarding

523 recognition results.

524 A comparison between proposed protocols and individual techniques high-
525 lighted that generating new samples by applying a selection of the best in-
526 dividual techniques, in a sequential one-by-one pipeline, is more convenient
527 than randomly picking one of them.

528 4.3. Comparison against existing methods

529 Results presented in Table 6 and Table 7 exhibit a considerable improve-
530 ment of the proposed method against the CBIA and ResNet methods in
531 terms of recognition performance. The results obtained by the ResNet are
532 poor in this context. This may be mainly due to the fact that the model
533 was originally intended to process images, and it lacks capabilities to learn
534 from temporal sequences as needed for this particular problem. It is impor-
535 tant to note here that results reported by Chelotti et al. (2018) are affected
536 by the use of a different tool to compare ground truth values against model
537 predictions. In this case, the temporal alignment of both events (real and
538 predicted) is considered using a gap or collar. By doing this, for example, a
539 sequence of events predicted in the correct order is not considered successful
540 if the temporal localisation does not match. Consequently, it is possible to
541 state that the comparison method proposed in this study is more rigorous
542 and appropriate for problems of JM event detection and classification.

543 In terms of computational costs, the proposed method involves a total of
544 464,919,007 floating point operations (FLOPs) in order to analyse one second
545 of the signal. The details about estimation of these costs are presented in the
546 Appendix A. This number represents an increase in the calculations needed
547 against the CBIA (1.000:1), which needs 398,860 FLOPs to process one sec-

548 ond of the signal. This value was estimated using the calculations reported
549 by the authors for the version (Least Mean Squares filter and Multi-Layer
550 Perceptron) and sampling frequency (22.05 kHz) used in the implementation
551 conducted here. Although the proposed method represents an increase in the
552 number of operations, the improvements obtained with respect to more ac-
553 curate recognition results represent a considerable advantage in the context
554 of applications where real-time operation is not required. The key advantage
555 of the proposed method is its ability to accurately classify JM using raw au-
556 dio signals, without any previous definition of sound features to be analysed
557 by the system. In this stage, the computational cost of algorithms is not
558 relevant compared with their ability to extract the appropriate information
559 without an “expensive”, handcrafted and generally non-optimal feature engi-
560 neering stage. This fact implies that this type of model can be used in the
561 development stages of a system when relevant features for JM recognition of
562 the sound are explored.

563 The interpretability of a proposed solution is another subject that must be
564 analysed from a practical point of view. In this sense, the method presented
565 in this paper poses a disadvantage when compared to traditional methods
566 that use "white box" models.

567 On the other hand, when algorithms must be deployed on IoT systems,
568 computational cost is a central issue since they must minimise the use of
569 energy. This type of operational condition requires that algorithms must be
570 optimised from the processor’s perspective, minimising the amount of energy
571 and memory as well as the notation used to represent the information. In
572 this way, handcrafted feature algorithms might require less implementation

573 effort in these scenarios. The price paid is the time and work required to
574 develop the system.

575 Concerning other DL methods, Li et al. (2021c) reported 88.8, 88.9 and
576 88.8 for F1 score, precision and recall respectively. Even though these values
577 seems to overcome the proposed Deep sound architecture in the classification
578 task, detection is disregarded in that study. Moreover, the limitations of the
579 approach proposed by Li et al. (2021c), plus the evaluation metrics proposed
580 here, should be considered in order to perform a direct comparison between
581 both methods. Finally, it is important to note that results reported by Li
582 et al. (2021c) slightly outperformed or was comparable to CBIA.

583 5. Conclusions

584 In this study, a novel end-to-end architecture for detection and classifica-
585 tion of ruminant masticatory JM events was presented and evaluated with
586 real data. The model combines two well known neural network types into a
587 single model, generating a CNN-RNN final architecture. Different numbers
588 of convolutional layers in the CNN block of the network have been explored.
589 The highest recognition performance (micro F1 score up to 79.8%) was ob-
590 tained using 4 pairs of convolution (plus dropout) layers. The use of data
591 augmentation has been evaluated, which resulted in an improvement of recog-
592 nition performance (almost 2.5% in terms of micro F1 score) when using a
593 selected subset of techniques to generate synthetic samples. The proposed
594 architecture outperformed a previous method (CBIA) by at least 10% (micro
595 F1 score) and a ResNet implementation by more than 30% (micro F1 score).
596 On the other hand, the proposed architecture automatically extracts features

597 from raw signals, which introduces very promising results when compared to
598 traditional methods that use manually created characteristics.

599 Future research will focus on the optimization of computational cost of
600 the proposed method, and the analysis of its impact on recognition results.
601 The interpretation of learned features and their corresponding qualitative
602 analysis will be part of future works. Finally, an exploration of transfer
603 learning, semi-supervised learning and related approaches will be studied in
604 order to evaluate other alternatives for small quantities of labelled data.

605 **Acknowledgments**

606 This work has been funded by Universidad Nacional del Litoral, CAID
607 50620190100080LI and 50620190100151LI, Universidad Nacional de Rosario,
608 projects 2013-AGR216, 2016-AGR266 and 80020180300053UR, Agencia San-
609 tafesina de Ciencia, Tecnología e Innovación (ASACTEI), project IO-2018-
610 -00082, Consejo Nacional de Investigaciones Científicas y Técnicas (CON-
611 ICET), project 2017-PUE sinc(i). Authors would like to thank the dedication
612 and perceptive help by Campo Experimental J. Villarino Dairy Farm staff
613 for their assistance and support during the completion of this study. Au-
614 thors also gratefully acknowledge the support of NVIDIA Corporation with
615 the donation of the Titan XP GPU used for this research.

616 **References**

- 617 Andriamandroso, A., Bindelle, J., Mercatoris, B., and Lebeau, F. (2016). A
618 review on the use of sensors to monitor cattle jaw movements and behav-
619 ior when grazing. *Biotechnologie, Agronomie, Société et Environnement*,
620 20:273–286.
- 621 Andriamandroso, A., Lebeau, F., and Bindelle, J. (2015). Changes in biting
622 characteristics recorded using the inertial measurement unit of a smart-
623 phone reflect differences in sward attributes. In *7th Conference on Preci-
624 sion Livestock Farming*, pages 283–289.
- 625 Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S.,
626 Barbado, A., García, S., Gil-López, S., Molina, D., Benjamins, R., et al.
627 (2020). Explainable artificial intelligence (xai): Concepts, taxonomies,
628 opportunities and challenges toward responsible ai. *Information fusion*,
629 58:82–115.
- 630 Bahmei, B., Birmingham, E., and Arzanpour, S. (2022). CNN-RNN and
631 data augmentation using deep convolutional generative adversarial net-
632 work for environmental sound classification. *IEEE Signal Processing Let-
633 ters*, 29:682–686.
- 634 Balasso, P., Marchesini, G., Ughelini, N., Serva, L., and Andrighetto, I.
635 (2021). Machine learning to detect posture and behavior in dairy cows:
636 Information from an accelerometer on the animal’s left flank. *Animals*,
637 11(10):2972.

- 638 Balch, C. (1958). Observations on the act of eating in cattle. *British Journal*
639 *of Nutrition*, 12(3):330–345.
- 640 Calamari, L., Soriani, N., Panella, G., Petrera, F., Minuti, A., and Trevisi,
641 E. (2014). Rumination time around calving: An early signal to detect cows
642 at greater risk of disease. *Journal of Dairy Science*, 97(6):3635–3647.
- 643 Chelotti, J. O., Vanrell, S. R., Galli, J. R., Giovanini, L. L., and Rufiner, H. L.
644 (2018). A pattern recognition approach for detecting and classifying jaw
645 movements in grazing cattle. *Computers and Electronics in Agriculture*,
646 145:83–91.
- 647 Chelotti, J. O., Vanrell, S. R., Milone, D. H., Utsumi, S. A., Galli, J. R.,
648 Rufiner, H. L., and Giovanini, L. L. (2016). A real-time algorithm for
649 acoustic monitoring of ingestive behavior of grazing cattle. *Computers*
650 *and Electronics in Agriculture*, 127:64–75.
- 651 Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F.,
652 Schwenk, H., and Bengio, Y. (2014). Learning phrase representations using
653 RNN Encoder-Decoder for statistical machine translation. In *Proceedings*
654 *of the Empirical Methods in Natural Language Processing (EMNLP 2014)*.
655 arXiv.
- 656 De Boever, J., Andries, J., De Brabander, D., Cottyn, B., and Buysse, F.
657 (1990). Chewing activity of ruminants as a measure of physical struc-
658 ture—a review of factors affecting it. *Animal Feed Science and Technology*,
659 27(4):281–291.

- 660 Ding, L., Lv, Y., Jiang, R., Zhao, W., Li, Q., Yang, B., Yu, L., Ma, W., Gao,
661 R., and Yu, Q. (2022). Predicting the feed intake of cattle based on jaw
662 movement using a triaxial accelerometer. *Agriculture*, 12(7):899.
- 663 Fogarty, E. S., Swain, D. L., Cronin, G. M., Moraes, L. E., and Trotter, M.
664 (2020). Behaviour classification of extensively grazed sheep using machine
665 learning. *Computers and Electronics in Agriculture*, 169:105175.
- 666 Frost, A. R., Schofield, C. P., Beulah, S. A., Mottram, T. T., Lines, J. A.,
667 and Wathes, C. M. (1997). A review of livestock monitoring and the
668 need for integrated systems. *Computers and Electronics in Agriculture*,
669 17(2):139–159.
- 670 Giovanetti, V., Decandia, M., Molle, G., Acciaro, M., Mameli, M., Cabiddu,
671 A., Cossu, R., Serra, M., Manca, C., Rassu, S., et al. (2017). Automatic
672 classification system for grazing, ruminating and resting behaviour of dairy
673 sheep using a tri-axial accelerometer. *Livestock Science*, 196:42–48.
- 674 Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training
675 deep feedforward neural networks. In *Proceedings of the thirteenth inter-*
676 *national conference on artificial intelligence and statistics*, pages 249–256.
677 JMLR Workshop and Conference Proceedings.
- 678 Hershey, S., Chaudhuri, S., Ellis, D. P., Gemmeke, J. F., Jansen, A., Moore,
679 R. C., Plakal, M., Platt, D., Saurous, R. A., Seybold, B., et al. (2017). CNN
680 architectures for large-scale audio classification. In *2017 ieee international*
681 *conference on acoustics, speech and signal processing (icassp)*, pages 131–
682 135. IEEE.

- 683 Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. R. (2012). Improving neural networks by preventing co-adaptation
684 of feature detectors.
685
- 686 Hoxhallari, K; Purcell, W. N. T. (2022). Precision livestock farming. In *10th*
687 *European Conference on Precision Livestock Farming*.
- 688 Kamminga, J. W., Le, D. V., Meijers, J. P., Bisby, H., Meratnia, N., and
689 Havinga, P. J. (2018). Robust sensor-orientation-independent feature selection for animal activity recognition on collar tags. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 2(1):1–27.
690
691
- 692 Khamees, A. A., Hejazi, H. D., Alshurideh, M., and Salloum, S. A. (2021).
693 Classifying audio music genres using CNN and RNN. In Hassanien, A.-
694 E., Chang, K.-C., and Mincong, T., editors, *Advanced Machine Learning Technologies and Applications*, pages 315–323, Cham. Springer International Publishing.
695
696
- 697 Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization.
698
- 699 Kokalis, C.-C. A., Tasakos, T., Kontargyri, V. T., Siolas, G., and Gonos, I. F. (2020). Hydrophobicity classification of composite insulators based
700 on convolutional neural networks. *Engineering Applications of Artificial Intelligence*, 91:103613.
701
702
- 703 Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2017). Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90.
704
705

- 706 Lecun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-
707 based learning applied to document recognition. *Proceedings of the IEEE*,
708 86(11):2278–2324.
- 709 Li, C., Tokgoz, K. K., Fukawa, M., Bartels, J., Ohashi, T., Takeda, K.-i.,
710 and Ito, H. (2021a). Data augmentation for inertial sensor data in cnns
711 for cattle behavior classification. *IEEE Sensors Letters*, 5(11):1–4.
- 712 Li, D., Liu, J., Yang, Z., Sun, L., and Wang, Z. (2021b). Speech emotion
713 recognition using recurrent neural networks with directional self-attention.
714 *Expert Systems with Applications*, 173:114683.
- 715 Li, G., Xiong, Y., Du, Q., Shi, Z., and Gates, R. S. (2021c). Classifying
716 ingestive behavior of dairy cows via automatic sound recognition. *Sensors*,
717 21(15).
- 718 Lim, S. J., Jang, S. J., Lim, J. Y., and Ko, J. H. (2019). Classification of
719 snoring sound based on a recurrent neural network. *Expert Systems with*
720 *Applications*, 123:237–245.
- 721 Lu, R., Duan, Z., and Zhang, C. (2018). Multi-scale recurrent neural network
722 for sound event detection. In *2018 IEEE International Conference on*
723 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 131–135.
- 724 Martinez-Rau, L. S., Chelotti, J. O., Vanrell, S. R., Galli, J. R., Utsumi,
725 S. A., Planisich, A. M., Rufiner, H. L., and Giovanini, L. L. (2022). A
726 robust computational approach for jaw movement detection and classifi-
727 cation in grazing cattle using acoustic signals. *Computers and Electronics*
728 *in Agriculture*, 192:106569.

- 729 Matsui, K. and Okubo, T. (1991). A method for quantification of jaw move-
730 ments suitable for use on free-ranging cattle. *Applied Animal Behaviour*
731 *Science*, 32(2-3):107–116.
- 732 Meng, J., Wang, X., Wang, J., Teng, X., and Xu, Y. (2022). A capsule
733 network with pixel-based attention and BGRU for sound event detection.
734 *Digital Signal Processing*, 123:103434.
- 735 Mesaros, A., Heittola, T., Virtanen, T., and Plumbley, M. D. (2021). Sound
736 event detection: A tutorial. *IEEE Signal Processing Magazine*, 38(5):67–
737 83.
- 738 Milone, D. H., Galli, J. R., Cangiano, C. A., Rufiner, H. L., and Laca, E. A.
739 (2012). Automatic recognition of ingestive sounds of cattle based on hidden
740 markov models. *Computers and Electronics in Agriculture*, 87:51–55.
- 741 Milone, D. H., Rufiner, H. L., Galli, J. R., Laca, E. A., and Cangiano,
742 C. A. (2009). Computational method for segmentation and classification
743 of ingestive sounds in sheep. *Computers and Electronics in Agriculture*,
744 65(2):228–237.
- 745 Monteiro, A., Santos, S., and Gonçalves, P. (2021). Precision agriculture for
746 crop and livestock farming—brief review. *Animals*, 11(8):2345.
- 747 Nanni, L., Paci, M., Brahnam, S., and Lumini, A. (2021). Comparison
748 of different image data augmentation approaches. *Journal of Imaging*,
749 7(12):254.
- 750 Navon, S., Mizrach, A., Hetzroni, A., and Ungar, E. D. (2013). Automatic

- 751 recognition of jaw movements in free-ranging cattle, goats and sheep, using
752 acoustic monitoring. *Biosystems Engineering*, 114(4):474–483.
- 753 Neethirajan, S. (2020). The role of sensors, big data and machine learning
754 in modern animal farming. *Sensing and Bio-Sensing Research*, 29:100367.
- 755 Nydegger, F., Gyga, L., and Egli, W. (2011). Automatic measurement of jaw
756 movements in ruminants by means of a pressure sensor. In *International
757 Conference on Agricultural Engineering*, page 27.
- 758 Oudshoorn, F. W., Cornou, C., Hellwing, A. L. F., Hansen, H. H., Munks-
759 gaard, L., Lund, P., and Kristensen, T. (2013). Estimation of grass intake
760 on pasture for dairy cows using tightly and loosely mounted di- and tri-
761 axial accelerometers combined with bite count. *Computers and Electronics
762 in Agriculture*, 99:227–235.
- 763 Papakipos, Z. and Bitton, J. (2022). Augly: Data augmentations for robust-
764 ness. *arXiv preprint arXiv:2201.06494*.
- 765 Paudyal, S., Maunsell, F. P., Richeson, J. T., Risco, C. A., Donovan, D. A.,
766 and Pinedo, P. J. (2018). Rumination time and monitoring of health dis-
767 orders during early lactation. *Animal*, 12(7):1484–1492.
- 768 Penning, P. D. (1983). A technique to record automatically some aspects
769 of grazing and ruminating behaviour in sheep. *Grass and Forage Science*,
770 38(2):89–96.
- 771 Petmezas, G., Cheimariotis, G.-A., Stefanopoulos, L., Rocha, B., Paiva,
772 R. P., Katsaggelos, A. K., and Maglaveras, N. (2022). Automated lung

- 773 sound classification using a hybrid CNN-LSTM network and focal loss
774 function. *Sensors*, 22(3):1232.
- 775 Ramirez, A. E., Donati, E., and Chousidis, C. (2022). A siren identification
776 system using deep learning to aid hearing-impaired people. *Engineering*
777 *Applications of Artificial Intelligence*, 114:105000.
- 778 Riaboff, L., Shalloo, L., Smeaton, A., Couvreur, S., Madouasse, A., and
779 Keane, M. (2022). Predicting livestock behaviour using accelerometers: A
780 systematic review of processing techniques for ruminant behaviour predic-
781 tion from raw accelerometer data. *Computers and Electronics in Agricul-*
782 *ture*, 192:106610.
- 783 Rombach, M., Südekum, K.-H., Münger, A., and Schori, F. (2019). Herbage
784 dry matter intake estimation of grazing dairy cows based on animal, be-
785 havioral, environmental, and feed variables. *Journal of Dairy Science*,
786 102(4):2985–2999.
- 787 Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning
788 representations by back-propagating errors. *Nature*, 323(6088):533–536.
- 789 Ruuska, S., Kajava, S., Mughal, M., Zehner, N., and Mononen, J. (2016).
790 Validation of a pressure sensor-based system for measuring eating, rumi-
791 nation and drinking behaviour of dairy cattle. *Applied Animal Behaviour*
792 *Science*, 174:19–23.
- 793 Schuster, M. and Paliwal, K. K. (1997). Bidirectional recurrent neural net-
794 works. *IEEE Transactions on Signal Processing*, 45(11):2673–2681.

- 795 Shorten, C. and Khoshgoftaar, T. M. (2019). A survey on image data aug-
796 mentation for deep learning. *Journal of Big Data*, 6(1).
- 797 Simonyan, K. and Zisserman, A. (2014). Very deep convolutional networks
798 for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- 799 Sokolova, M. and Lapalme, G. (2009). A systematic analysis of performance
800 measures for classification tasks. *Information Processing and Management*,
801 45(4):427–437.
- 802 Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. (2016).
803 Rethinking the inception architecture for computer vision. In *Proceedings*
804 *of the IEEE conference on computer vision and pattern recognition*, pages
805 2818–2826.
- 806 Tani, Y., Yokota, Y., Yayota, M., and Ohtani, S. (2013). Automatic recogni-
807 tion and classification of cattle chewing activity by an acoustic monitoring
808 method with a single-axis acceleration sensor. *Computers and Electronics*
809 *in Agriculture*, 92:54–65.
- 810 Ungar, E. D., Ravid, N., Zada, T., Ben-Moshe, E., Yonatan, R., Baram, H.,
811 and Genizi, A. (2006). The implications of compound chew–bite jaw move-
812 ments for bite rate in grazing cattle. *Applied Animal Behaviour Science*,
813 98(3-4):183–195.
- 814 Vanrell, S. R., Chelotti, J. O., Bugnon, L. A., Rufiner, H. L., Milone, D. H.,
815 Laca, E. A., and Galli, J. R. (2020). Audio recordings dataset of grazing
816 jaw movements in dairy cattle. *Data Brief*, 30:105623.

817 Werner, J., Leso, L., Umstatter, C., Niederhauser, J., Kennedy, E., Geoghe-
818 gan, A., Shalloo, L., Schick, M., and O'Brien, B. (2018). Evaluation of
819 the rumiwatchsystem for measuring grazing behaviour of cows. *Journal of*
820 *Neuroscience Methods*, 300:138–146.

821 Wilcoxon, F. (1945). Individual comparisons by ranking methods. *Biometrics*
822 *Bulletin*, 1(6):80.

823 Zhu, Z., Dai, W., Hu, Y., and Li, J. (2020). Speech emotion recognition model
824 based on bi-gru and focal loss. *Pattern Recognition Letters*, 140:358–365.

825 **Appendix A. Computational costs**

826 The amount of operations required for processing one second of audio
827 signal were estimated at a sampling frequency of 6 kHz, a time window of
828 300 ms and a hop length of 150 ms. The procedure used to estimate these
829 calculations is similar to the one used in Chelotti et al. (2018) in which
830 additions and multiplications count as separated operations. The model
831 architecture presented in Figure 4 (c) was used here for comparison purposes.

832 In the first block of the proposed model, the following layers were con-
833 sidered: re-scaling, 1D convolution and max pooling. FLOPs required for
834 activation functions were also considered. Dropouts were discarded because
835 these layers only applied during training, and no calculations were considered
836 for the flatten operation. The cost of each of the convolutional layers were
837 estimated using the following expression:

$$(2 * C_i * K * H * W * C_o) \tag{A.1}$$

838 where C_i and C_o represents the input and output channels, K the kernel
839 size, H and W the size of the output feature map. According to this, the
840 total number of FLOPs in the first block of the model is 272.235.413.

841 In the second block of the model, FLOPs involved in reset and update
842 gates, activation functions and output generation were considered for every
843 unit. The total number of FLOPs required is 191.363.413.

844 Finally, in the last block of the model, the FLOPs required in dense layers
845 as well as activation functions were considered. The cost of each dense layer
846 were estimated using the following expression:

$$(2 * I * O) \tag{A.2}$$

847 where I and O represent the number of input and output neurons, re-
848 spectively. The total number of FLOPs in the last block of the model is
849 1.320.180. In summary, the total number of FLOPs in order to process one
850 second of signal is 464.919.007.