

Noise, Novels, Numbers. A Framework for Detecting and Categorizing Noise in Danish and Norwegian Literature

Anonymous ACL submission

Abstract

We present a framework for detecting and categorizing noise in literary texts, demonstrated through its application to Danish and Norwegian literature from the late 19th century. Noise, understood as “aberrant sonic behaviour,” is not only an auditory phenomenon but also a cultural construct tied to the processes of civilization and urbanization. By leveraging topic modeling techniques and fine-tuned BERT-based language models trained on Danish and Norwegian texts, we analyze a corpus of over 800 novels to extract and examine noise-related topics. We identify and track the prevalence of noise in these texts, offering insights into the literary perceptions of noise during the Scandinavian “Modern Breakthrough” period (1870-1899). Our contributions include the development of a comprehensive dataset annotated for noise-related segments and their categorization into human-made, non-human-made, and musical noises. This study illustrates the framework’s potential for enhancing the understanding of the relationship between noise and its literary representations, providing a deeper appreciation of the auditory elements in literary works.

1 Introduction

Noise, understood as “deviant sonic behaviour” (Novak and Sakakeeny, 2015), is an auditory phenomenon, but also a cultural construct, closely tied to processes of civilization and urbanization. The representation of noise in literature provides insight into the social and cultural developments of the period(s) in which that literature was written.

But can we trace how the past may have sounded? In studies of literature and sound, the empirical foundation is often a small selection of texts, representing either canonical and/or avant-garde instances of 20th century modernism (e.g., Toth, 2016; Lambrecht, 2017; Frattarola, 2018). In our investigations of the soundscapes in the litera-

ture of the Scandinavian “Modern Breakthrough” (1870-1899), we broaden the empirical and cultural scope by reading at scale. For that, we develop a framework for the detection and categorization of noise in literary texts.

We extract a series of noise-related topics from a corpus containing more than 800 Danish and Norwegian novels. By examining changes in the frequency of these topics, we draw insight into how literary perceptions of noise have evolved. The findings of this study contribute to a deeper understanding of the relationship between noise and its various representations in literary contexts.

Contributions. Our contributions are: (1) the development of a robust and scalable framework for detecting and categorizing noise in literary texts using fine-tuned language models and topic models; (2) the creation of an annotated dataset derived from over 800 Danish and Norwegian novels from the late 19th century, with detailed classifications of noise-related segments into human-made, non-human made, and musical noises; (3) the implementation of analysis to track the evolution of noise perceptions in literature over time, providing valuable insights into the cultural and social changes reflected in the soundscapes of the Scandinavian ‘Modern Breakthrough’; and (4) a demonstration of the framework’s applicability and effectiveness in literary studies, paving the way for future research in other linguistic and cultural contexts.

2 Related Work

Topic modeling on diachronic text. Various methods have been explored for topic modeling in digital humanities and computational literary analysis, specifically in diachronic analysis. Challenges in analyzing diachronic data using topic models are discussed in Marjanen et al. (2020), which presents a method for applying topic models to large and imbalanced collections. Heyer et al.

(2017) explore detecting semantic change and term evolution through three approaches and introduce “context volatility” as a measure for detecting semantic change. The “Draw My Topics” toolkit, presented by Dou et al. (2016), uses an algorithm based on Vector Space Model and Conditional Entropy to incorporate social scientists’ interests into standard topic modeling. Haider (2019) use Latent Dirichlet Allocation (LDA) for distant reading tasks on literary data, classifying poems into time periods and authorship attribution, while Tangherlini and Leonard (2013) show how LDA can be used as part of a search and discovery pipeline for literary study and the emergence of topics across domains. Finally, Roberts et al. (2013)’s Structural Topic Models use metadata to generate topic prevalence and have had a substantial impact in the social sciences.

Noise/sound in literary texts. A strand of inquiry close to our empirical approach focuses on the representation of sounds and their associated soundscapes within historical and fictional worlds. Schafer (1993) pioneers this approach by examining evolving soundscapes, both real-world and fictional, emphasizing the role of writers as “ear-witnesses” to past sonic environments. Recent studies extend this exploration to fictional soundscapes, exploring how descriptions of ambient sounds contribute to immersive storytelling experiences (Verma, 2019; Mildorf, 2019).

Further investigations explore specific genres, such as Gothic fiction, to uncover how ambient soundscapes shape narrative atmospheres and reader perceptions (Guhr and Algee-Hewitt, 2024). Fine-tuned BERT models successfully detect and analyze word-level sound indicators across literary texts. This interdisciplinary approach not only enhances our understanding of sound’s role in genre classification but also sheds light on broader narrative techniques and reader engagement strategies.

3 Scandinavian Literary Soundscapes

Studies in sensory history, and particularly *sound-scapes* studies, stress the importance of the sonic environments in which people live and operate. The importance of these soundscapes is based on the premise that the sounds one hears in a given place are as distinctive and as important as the things one sees there (Birdsall, 2012). 19th century urbanized settings presaged a ‘sound revolution,’ where novel individual sounds and soundscapes

rapidly emerged due to new industries, machinery, means of transportation, road surfaces, and the like (Parby, 2021). Simultaneously a fundamental shift in people’s relationship with sound and noise took place, leading some scholars to characterize the 19th and early 20th centuries as a particular auscultative age—an era devoted to sonic experiences, to ‘close listening’ and to auscultation “not only in the medical sense initiated by the stethoscope [...] and perfected by the microphone [...] but also in the sense of careful listening to a world at large.” (Picker, 2003).

In a Scandinavian context, these sensory events are related to the “modern breakthrough.” Copenhagen changed radically during this time, with an explosion in its population, and the introduction of new technologies and infrastructures. As Parby (2021) has shown, these developments led to the emergence of new soundscapes, which were then incorporated into contemporary fiction. Simultaneously, authors developed new realist literary techniques to give a fuller account of the material world and to represent it with verisimilitude and detail (Bjerring-Hansen and Wilkens, 2023).

4 Defining Noise

In Sound Studies, the definition and phenomenological demarcation of noise has been a point of considerable discussion. Etymologically, the word ‘noise’ is rooted in Latin ‘nausea’ that encompasses seasickness and nausea, and refers to sounds that are perceived as excessive, incoherent, confused, or twisted. One of the founders of Sound Studies, Raymond Murray Schafer, proposed a clear dichotomy between natural and man-made noises from, for instance, industrial activity and traffic (Schafer, 1993). This position has later been criticized for being a far too normative and rigid division, where man-made noise almost by definition is perceived as negative and natural sound events such as thunder claps as positive (Kelman, 2010).

Some more subject-sensitive, less normative and, not least, operational definitions, which we rely on, have been suggested by David Novak who defines noise as “deviant sonic behaviour” (Novak and Sakakeeny, 2015), and Peter Bailey who defines it as ‘sound out of place’ paraphrasing anthropologist Mary Douglas’ classic definition of dirt as ‘matter out of place’ (Bailey, 1996).

In order to develop an operational conceptual framework, we use a very basic definition of noise

as “silence-breaking”; this characterization means that noise can include moderate sounds such as whispering or mumbling voices as long as they are noted as sonic events in the text. It also means that we include neutral sonic phenomena, such as this description of factory whistles:

The steam pipes sounded from all the factories, it was 8 o'clock. And I had to leave. (transl.)

The quote says nothing about whether the whistles are perceived as noise by the character or narrator. But, obviously, the sound distinctively breaks a silence.

Based on our knowledge of industrialization and urbanization, as discussed in international Sound Studies research and reflected in literary sources, our hypothesis in the following is that during the latter part 19th century the general noise levels increased, correlated with an increasing sensitivity towards noise, in a development more and more dominated by non-human noise sources.

5 An Annotated Dataset of Noise in 19th Century Scandinavian Literature

We introduce a framework for noise detection and categorization, and apply it to a corpus of Scandinavian literature, creating a new annotated dataset.

5.1 Main Corpus

For our main target data, we rely on the MeMo corpus (Bjerring-Hansen et al., 2022), comprising 859 Danish and Norwegian novels spanning the last 30 years of the 19th century, with more than 64 million tokens. We refer to this corpus as the ‘main corpus’. The corpus is a rich and diverse collection of texts that provides valuable insights into the representations of noise and sound during the period under investigation. Table 1 shows statistical information about the corpus. We segment the corpus into paragraphs and split them into 50-word segments if they exceed 50 words.

Total novels	859
Total segments	1,936,527
Total words	64,227,927
Average segments per novel	2,254
Average words per novel	74,771
Average words per segment	33

Table 1: MeMo corpus statistics.

5.2 Noise Detection Dataset

To construct a dataset of text segments annotated for whether they contain noise or not, we combine a selection of hand-picked segments by a historian expert, with a topic-based search approach to enrich the dataset.

Segments Extraction. We apply BERTopic (Grootendorst, 2022), a powerful topic modeling technique that enables us to cluster millions of text segments from the MeMo corpus into a concise set of topics. By doing so, we aim to distill vast amounts of textual data into manageable thematic clusters, facilitating subsequent analysis. Following the topic modeling phase, we filter the generated topics, focusing our attention on those most relevant to the concept of noise. These selected topics serve as a foundation for further exploration, guiding us in identifying and annotating text segments specifically related to noise. We remain with 5,700 text segments potentially related to noise topics, which are then carefully annotated by experts.

Annotation Guidelines. The annotation involves two of the authors, native Danish speakers, a historian with special interests in urban and social history as well as a literary scholar familiar with the social conditions of 19th century literature, who classified segments into two categories: noise and non-noise. This decision is by no means trivial.

With the conceptual distinctions and demarcations in §4 in mind, the annotation of noise and non-noise text segments is carried out on the basis of the following, minimalist and pragmatic guidelines, respecting the principle that clear and simple instructions are crucial for obtaining high-quality annotations (Mohammad, 2016), while also acknowledging the intricacies that the analysis of literary texts based on small fragmentary segments raises.

1. Based on our definition of noise, the text segments are labeled either ‘1’ (positive) or ‘0’ (negative). Our focal point is that we would rather narrow the focus later on (through additional annotation of positive cases, metadata-filtering, or NLP measures) than exclude specific types of noise. Along the way, our definition become even richer, as we realize that music should also be included, as music is often interwoven with other types of sonic events, as we see in this example:

The orchestra began playing a French folk tune ... and the stormy applause and

269	applause clapping eventually faded away.
270	(transl.)
271	2. Only the segment in question should be con-
272	sidered. Contextualisation and ‘guessing’ on
273	what might go on before or after the segment
274	were ruled out. In cases of doubt, the label
275	should be ‘0’. Example:
276	But after the noisy scenes he had caused,
277	came the lethargy that always follows the
278	performance of a great tragic role. (transl.)
279	3. Negated noise should be filtered out, e.g,
280	The gardens alongside the houses looked
281	very well on this summer evening. There
282	was no noise or disturbance, only a couple
283	of children playing across the street, but
284	they did not chatter; even their play was
285	in keeping with the tone of the evening.
286	(transl.)
287	4. The same goes for other pseudo-relevant seg-
288	ments detected by the topic modeling algo-
289	rithm, including noise metaphors and similies
290	(prompted by words like ‘as’ or ‘like’), as we
291	see in this example:
292	She felt engulfed by a buzzing electric cur-
293	rent. (transl.)
294	Annotation Results. Our hand-picked selection
295	of noise segments include 217 positive examples,
296	manually curated from various 19 th century sources
297	(memoirs, essays and fictional works). As for the
298	segments obtained with the topic-based search, out
299	of the pool of 5700 segments, 337 are deemed
300	noise-related while the remaining segments are an-
301	notated as non-noise.
302	In our endeavor to encompass the entirety of the
303	target corpus, we turn our attention to the remaining
304	5,365 segments, considering them as “non-noise”
305	segments. To ensure the dataset’s comprehensiveness
306	and diversity, we supplement these non-noise
307	segments with an additional 5,000 segments ran-
308	domly selected from various non-noise topics.
309	We randomly sample 175 segments from the
310	noise-related topics to serve as our testing set, each
311	annotated independently by our two annotators, to
312	evaluate annotation consistency and assess model
313	performance. The resulting Cohen’s Kappa value,
314	calculated to measure inter-annotator agreement,
315	yielded a score of 0.85, indicating a high level of
316	agreement between the annotations.

5.3 Noise Categorization Dataset	317
Fine-grained categorization of noise-related seg-	318
ments is essential in the context of classifying tex-	319
tual noise extracted from literary texts. This ap-	320
proach enables a nuanced understanding of the di-	321
verse forms of noise present within the textual cor-	322
pus, including but not limited to, linguistic anom-	323
alies, contextual inconsistencies, and stylistic irreg-	324
ularities. By classifying noise into categories such	325
as human noise, mechanical noise, and other types	326
of textual disturbances, one can distinguish specific	327
sources of interference more accurately, reflecting	328
the intricacies inherent in literary compositions.	329
Annotation Guidelines. To get closer to an un-	330
derstanding of sound as a cultural phenomenon	331
as reflected in literary works, we perform another	332
round of annotation. Although we have several	333
specific research interests related to sonic develop-	334
ments in the 19 th century, in order to (a) reduce the	335
number of axes in the annotation, which might have	336
negative consequences for the predictive power of	337
the model, and (b) produce a more broadly useful	338
dataset, we choose to prioritize one aspect of the	339
noise segments, namely the <i>sound source</i> . In do-	340
ing so, we have disregarded features, which, in a	341
future, extended pipeline, may come into play, not	342
least <i>time</i> (traditional or modern sound?) and <i>place</i>	343
(rural or urban sound?).	344
For this round of annotation, we merge the noise-	345
related segments from the previous dataset and ad-	346
ditional segments from the MeMo corpus after the	347
prediction of noise and non-noise classes for each	348
segment in the corpus as shown in Figure 1. Then,	349
the (same two) annotators classified text segments	350
into the following categories: Non-human made	351
noise (T), Human-made noise (H), Undefined noise	352
(N), and Music (M), following these criteria to en-	353
sure an accurate and consistent categorization:	354
1. Non-human made noise encompasses any	355
noise not produced by humans, ranging from	356
machine-produced sonic events (such as steam	357
engines, trams, telephones etc.) to natural	358
ones (caused by wind, rain, animals etc.).	359
2. Human-made noise includes any noise result-	360
ing from human activities (such as footsteps,	361
conversation, yelling, booing etc.).	362
3. Undefined noise is the appropriate label when	363
the noise source is unknown or unclear, as in	364

365 this example where the noise is abstract and
366 generic:

367 Outside, the music stopped and Madsen’s
368 voice was barely audible through the noise.
369 (transl.)

370 4. Music is a special category in relation to
371 both the general ontology: noise yes/no?
372 (see above) and to the specific categorization.
373 Since it is futile to determine whether a rat-
374 tling sound is produced by the violin player
375 or his instrument, we decided to give music a
376 label of its own.

377 5. Often there is a mix of sound sources in the
378 individual text segments, as here (**non-human**
379 **made**, **human made**, **music**):

380 From time to time there were **snatches of a**
381 **loud violin’s dance tunes**, **Lonely cabs rum-**
382 **bled** through the street, with **snow-damped**
383 **wheels and a few swishing whistles** that
384 made a couple of heads turn in the window.
385 Every now and then **a streetcar threw its**
386 **jingle of bells and chimes** into **the whispers**
387 **and murmurs of conversation**. (transl.)

388 Here, non-human made noise dominates, so
389 the label is ’T’. In other cases, the categoriza-
390 tion of mixed segments is based on a more
391 uncertain basis and is open to interpretation.

392 **Annotation Results.** Following the annotation
393 process, we have a total of 1,874 text segments
394 annotated by two independent annotators. Annot-
395 ated data statistics are presented in Table 2. The
396 training set, encompassing ~91% of the total anno-
397 tations, consists of 1,699 segments, while the test
398 set, comprising ~9%, consists of 175 segments.
399 After the removal of non-noise segments, the total
400 number of segments in the dataset is 1,244. No-
401 tably, both annotators annotate all segments within
402 the test set, ensuring comprehensive coverage and
403 reliability. Our obtained Cohen’s Kappa value of
404 ~0.81 demonstrates a substantial agreement level,
405 surpassing chance expectations. This result under-
406 scores the robust and accurate classification of the
407 data, reflecting strong and reliable consistency in
408 the annotations provided by both annotators. Note
409 that in the training set, each annotator individually
410 annotates half of the segments, maintaining an eq-
411 uitable distribution to uphold annotation quality
412 and consistency across the dataset.

Non-human	Human-made	Undefined	Music
513 (40%)	424 (33%)	56 (4%)	269 (21%)

Table 2: Noise categorization annotated data statistics.

6 Experiments and Results 413

414 In this section, we describe the selection of pre-
415 trained language models as well as the classifica-
416 tion experiments on the noise detection and noise
417 categorization datasets.

6.1 Pre-trained Language Models 418

419 In this subsection, we outline the models evaluated
420 in our noise detection and categorization classifica-
421 tion experiments using supervised fine-tuning meth-
422 ods. Importantly, all models are selected based on
423 their performance evaluated on Danish and Norwe-
424 gian literary benchmark datasets (Al-Laith et al.,
425 2024), the Scandinavian Embedding Benchmark¹
426 and ScandEval² (Nielsen, 2023), even though these
427 models had not been trained primarily on historical
428 Danish or Norwegian.

429 **DanskBERT.** DanskBERT³, a top-performing
430 Danish language model noted for its success on
431 the ScandEval benchmark (Snæbjarnarson et al.,
432 2023), is based on the XLM-RoBERTa architec-
433 ture and trained on the Danish Gigaword Corpus
434 (Strømberg-Derczynski et al., 2021). It features 24
435 layers, a hidden dimension of 1024, 16 attention
436 heads, and a subword vocabulary of 250,000. The
437 model was trained with a batch size of 2,000 for
438 500,000 steps on 16 V100 GPUs over two weeks.

439 **Danish Foundation Models sentence encoder.**
440 A sentence-transformers model (Enevoldsen et al.,
441 2023) based on the BERT architecture, featuring
442 24 layers, 16 attention heads, and a hidden size of
443 1024. It incorporates a dropout rate of 0.1 for atten-
444 tion probabilities and hidden states, using GELU
445 activation and supporting up to 512 position em-
446 beddings. With a vocabulary size of 50,000 tokens,
447 this model, referred to as DFM (Large), excels in
448 tasks such as Danish sentiment analysis and named
449 entity recognition.⁴

¹<https://kennethenevoldsen.github.io/scandinavian-embedding-benchmark/>

²<https://scandeval.com/>

³<https://huggingface.co/vesteinn/DanskBERT>

⁴<https://huggingface.co/KennethEnevoldsen/dfm-sentence-encoder-large-exp2-no-lang-align>

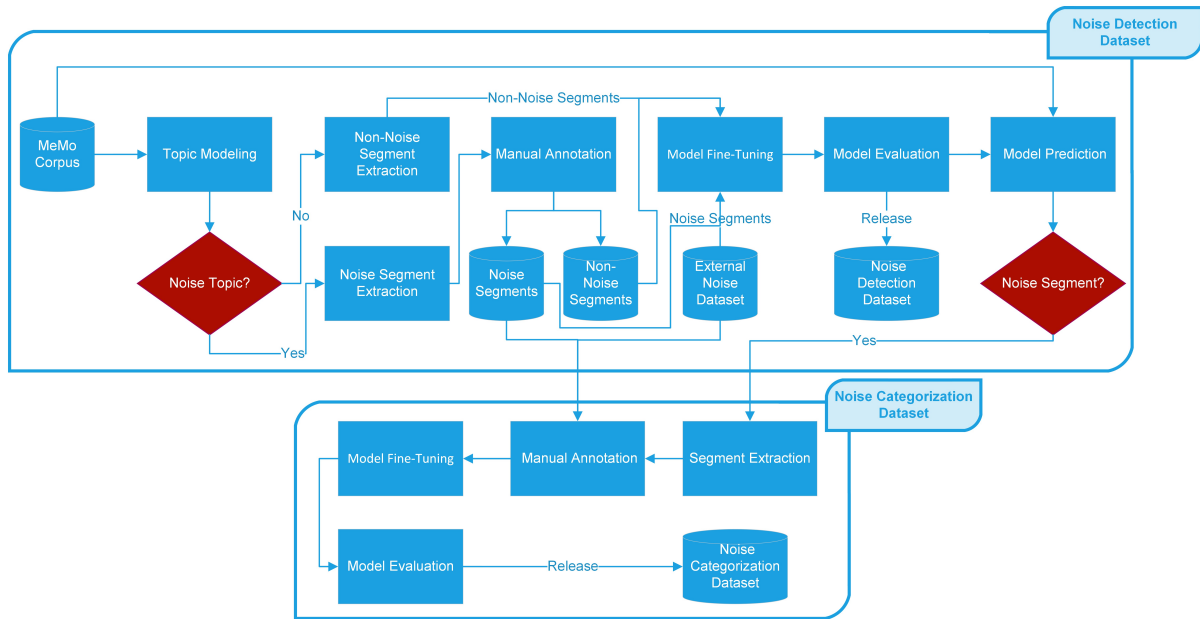


Figure 1: Noise Datasets Creation Flowchart.

MeMo-BERT-03. Developed by continuing the pre-training of the Transformer PLM DanskBERT (Al-Laith et al., 2024).⁵ This foundation allows MeMo-BERT-3 to leverage extensive linguistic knowledge for NLP tasks in historical literary Danish including sentiment analysis and word sense disambiguation. The model outperformed different models in sentiment analysis and word sense disambiguation tasks (Al-Laith et al., 2024).

NB-BERT-base. A general-purpose BERT-base model was developed using the extensive digital collection at the National Library of Norway (Kummervold et al., 2021).⁶ It follows the architecture of the BERT Cased multilingual model and has been trained on a diverse range of Norwegian texts, encompassing both Bokmål and Nynorsk from the past 200 years. This comprehensive training allows the NB-BERT-base to effectively handle a wide array of NLP tasks in Norwegian. The model achieved the second-highest performance ranking in the Norwegian Named Entity Recognition task compared to other models listed on the ScandEval benchmark for Norwegian natural language understanding.

6.2 Experimental Setup

In this section, we outline the experimental setup employed for the supervised classification tasks fo-

⁵<https://huggingface.co/MiMe-MeMo/MeMo-BERT-03>

⁶<https://huggingface.co/NbAiLab/nb-bert-base>

cused on both noise detection and noise categorization. Our experiments involve fine-tuning several pre-trained language models on the fine-grained datasets. The details of the dataset and the models used are described below. For the training procedure, the experiments involve fine-tuning BERT models on the dataset using a batch size of 32, training for 20 epochs with the AdamW optimizer at a learning rate of 10^{-3} . During training, we monitored both training and validation losses to assess model convergence and prevent overfitting. For evaluation, we employed the F1-score metric due to its ability to balance precision and recall, particularly effective for tasks with imbalanced datasets like noise detection and categorization. The performance of each model was evaluated on both validation and test sets, ensuring the robustness and generalizability of the models across different datasets and epochs.

6.3 Noise Detection Experiments

It is important to note the deliberate imbalance within the dataset, where only $\sim 5\%$ of the annotated segments are noise-related. By favoring a higher representation of non-noise segments, we aim to bias our model toward accurately identifying and capturing instances of noise within the data. This approach is designed to enhance the model’s sensitivity to noise while maintaining robustness in its classification capabilities.

Fine-tuning the PLMs on the noise detection task

507 results in notable performance variations (Table 3).
 508 DanskBERT achieves a validation accuracy of 0.89
 509 and a test accuracy of 0.83, indicating robust per-
 510 formance across unseen data. MeMo-BERT-03
 511 demonstrated the highest validation accuracy at
 512 0.90, although its test accuracy slightly decreased
 513 to 0.80. In contrast, DFM (Large) exhibited a vali-
 514 dation accuracy of 0.81, dropping significantly to
 515 0.55 on the test set, suggesting potential overfit-
 516 ting or limited generalizability. NB-BERT-base
 517 achieved consistent results with a validation accu-
 518 racy of 0.88 and a test accuracy of 0.76, indicating
 519 reliable performance across both validation and test
 520 datasets. These results highlight the effectiveness
 521 of fine-tuned BERT variants, especially MeMo-
 522 BERT-03 and DanskBERT, in accurately detecting
 523 noise within textual data, while emphasizing the
 524 importance of robust evaluation across multiple
 525 models and datasets.

526 6.4 Noise Categorization Experiments

527 The dataset comprises 1,244 text segments, divided
 528 into training, validation, and testing sets for model
 529 development and evaluation. The training set in-
 530 cludes 961 examples, constituting $\sim 77\%$ of the
 531 dataset, while the validation set, used for hyper-
 532 parameter selection, consists of 178 samples, rep-
 533 resenting $\sim 14\%$ of the total. The testing set, for
 534 the final model evaluation, contains 105 examples,
 535 or $\sim 9\%$ of the dataset. Annotations for the train-
 536 ing and validation sets were made by a single ex-
 537 pert. For the testing set, only segments where both
 538 experts agreed on the annotations were retained,
 539 discarding those with conflicting annotations. We
 540 use the weighted average F1-score as the evalua-
 541 tion metric. Notably, MeMo-BERT-03 achieved
 542 the highest F1-score of 83% on the validation set,
 543 while the DanskBERT model achieved the highest
 544 F1-score of 83% on the test set. Table 3 shows
 545 detailed results for each model.

546 7 Diachronic Analysis of Noise Segments

547 Having trained accurate noise detection and cate-
 548 gorization classifiers, we use the best ones (Dan-
 549 skBERT fine-tuned on the two tasks respectively) to
 550 predict labels for all segments in the entire MeMo
 551 corpus. We then quantify the frequency of the oc-
 552 currence of noise over time in the corpus, as well
 553 as the distribution of the different categories.

	Detection		Categorization	
Model	Valid.	Test	Valid.	Test
DanskBERT	0.89	0.83	0.81	0.83
DFM (Large)	0.81	0.55	0.82	0.79
MeMo-BERT-03	0.90	0.80	0.83	0.80
NB-BERT-base	0.88	0.76	0.81	0.81

Table 3: Validation and test F1-Score results of fine-tuning the selected models on the noise detection and categorization datasets.

554 7.1 Noise Occurrences over Time

555 After fine-tuning multiple pre-trained PLMs, Dan-
 556 skBERT emerged as the top performer among the
 557 four models evaluated, and we selected it for pre-
 558 dicting noise and non-noise classes across all seg-
 559 ments in the main corpus. Notably, out of 1.9
 560 million segments in the corpus, 22,0378 were pre-
 561 dicted with the noise class. Figure 2 shows the
 562 proportion of noise segments over the years.

563 A trend of rising noise levels in the novels is
 564 clear: From the 1870s to the 1890s there is a more
 565 than 50 % relative increase of noise (followed by a
 566 slight decline or a plateau by the end of the decade).

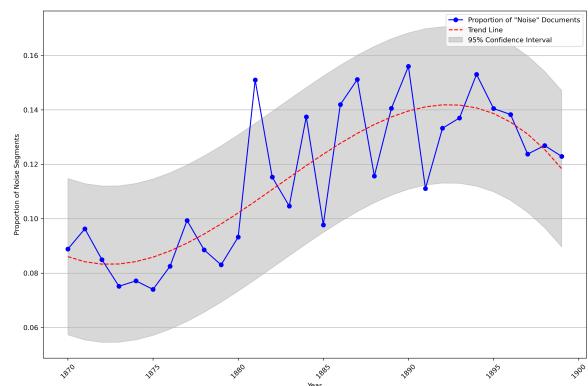


Figure 2: Proportion of Noise Segments Over Years.

567 7.2 Noise Categories over Time

568 Applying the best classifier for the second task
 569 (DanskBERT fine-tuned on noise categorization),
 570 we predict noise categories for all positive predic-
 571 tions in the corpus from the previous step. Figure 3
 572 shows the frequency of the noise categories.

573 The analysis indicates a stable distribution of the
 574 different kinds of noise without significant fluctua-
 575 tions.

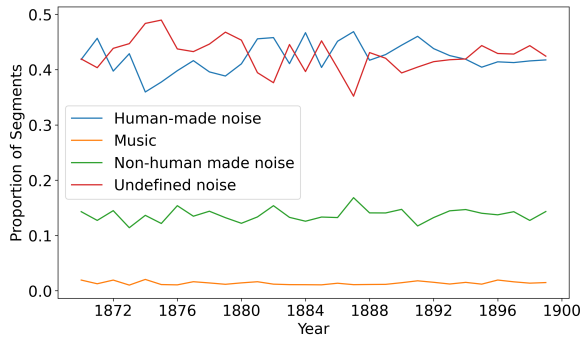


Figure 3: Frequency of Noise Categories over Time.

8 Discussion

Annotation nuances. Given the complex and slippery character of noise as a cultural phenomenon, we chose to disregard both spatial aspects (where?) and historical aspects (when?) related to it. For pragmatic reasons, we decided on a basic categorization by annotating noise classes (based on noise sources). There were challenges in making categorical decisions about specific noise events in the texts, for instance when there was a mix of noise sources simultaneously at play and, not least, in relation to our crucial distinction between non-human made and human made noise, according to which in the former category, it is humans themselves who generate sound (with their voice or body), while in the latter, technology is perceived as an agent and humans as mere operators (e.g. of a ringing church bell or a screeching streetcar). Our categories proved to be operational, but they are by no means watertight. We have learned more about noise as a historical concept, but also about its dynamic way of manifesting itself in different contexts and literary representations.

Model performance. Interestingly, while MeMoBERT-03 achieved the highest validation scores, DanskBERT outperformed it on the test set for both noise detection and categorization (Table 3), suggesting it is more capable at generalizing to unseen segments. While the former is designed to be better attuned to nuances in the historical corpus, in some cases, the latter (which was trained on a diverse corpus consisting mostly of modern Danish) might be better at detecting modernity signals, which are the focus of our annotation framework.

Noise trends. Our results confirm our hypothesis that in the last three decades of the 19th century a rise of general noise levels as well as an increase in preoccupation with noise are reflected in the novels

of the time. The upward trend is clear. It does, however, flatten or fall by the end of the period. This is hardly due to less noise, but rather to the fact that the authors and their characters are less preoccupied with it. Our hypothesis, which must be supported by close inspection and reading, is that noise is taken for granted or implied on the brink of the 20th century. In contrast to our initial hypothesis, human-made noise remains at a relatively high level throughout the period. This fits well with observations from the larger European metropolises like Paris, London and Madrid, where the policing of human noise sources remain significant, whereas industrial sounds tends to be evaluated positively and is not the focus of anti-noise campaigns until the early 20th century. The results from our noise categorization do however call for further time- and place-attentive investigation.

9 Conclusion

We presented a framework for detecting and categorizing noise in literary texts and demonstrated its usefulness in the MeMo corpus. Using topic modeling and fine-tuned BERT-based models, we extracted and analyzed relevant text segments, providing new insights into the cultural and social transformations reflected in the soundscapes of the Scandinavian “Modern Breakthrough” period. Our study demonstrates that literary perceptions of noise can be effectively tracked and categorized, revealing significant patterns and trends. We have been able to add new perspectives on the interplay between literature and cultural history – and to empirically underpin hypotheses of the 19th century as a particular auscultative era.

Future work will extend this framework to explore the impact of industrialization, examining how technological advancements and urbanization influenced literary soundscapes. We will also investigate the spatial dimensions of noise, contrasting rural and urban settings. Further, we would like to do comparative analysis on other datasets to situate our study in a broader context, such as contemporary civil complaints (as documented in the City Archives), and/or a corpus of modern Danish novels (from The World Literature Data Collective). Additionally, we plan to analyze the lexical diversity in the terms used to portray noise, to get a better understanding of the psychological and cognitive aspects of the increased awareness of noise in the early phases of urbanization and industrialization.

664 Limitations

665 Despite the strengths of our framework, there are
666 several limitations to consider. First, the focus on
667 Danish and Norwegian literature from a specific
668 historical period may limit the generalizability of
669 our findings to other linguistic and cultural con-
670 texts. Second, the accuracy of our noise detection
671 and categorization relies heavily on the quality of
672 the annotations and the pre-trained language mod-
673 els, which may not capture all nuances of noise
674 representation in literary texts. Third, our current
675 analysis does not account for the broader contextual
676 elements surrounding noise occurrences, such as
677 narrative structure or character perspectives, which
678 could provide a deeper understanding of the literary
679 soundscapes. Finally, while our framework demon-
680 strates promising results, further validation across
681 diverse datasets and more complex noise categor-
682 ization schemes is necessary to fully establish its
683 robustness and applicability.

684 Ethics Statement

685 The MeMo corpus, which we use, is released un-
686 der the Creative Commons Attribution 4.0 Interna-
687 tional license.

688 References

689 Ali Al-Laith, Alexander Conroy, Jens Bjerring-Hansen,
690 and Daniel Hershcovich. 2024. [Development and](#)
691 [evaluation of pre-trained language models for his-](#)
692 [torical Danish and Norwegian literary texts](#). In *Pro-*
693 *ceedings of the 2024 Joint International Conference*
694 *on Computational Linguistics, Language Resources*
695 *and Evaluation (LREC-COLING 2024)*, pages 4811–
696 4819, Torino, Italia. ELRA and ICCL.

697 Peter Bailey. 1996. Breaking the sound barrier: A histo-
698 rian listens to noise. *Body & Society*, 2(2):49–66.

699 Carolyn Birdsall. 2012. *Nazi soundscapes: sound, tech-*
700 *nology and urban space in Germany, 1933-1945*.
701 Amsterdam University Press.

702 Jens Bjerring-Hansen, Ross Deans Kristensen-
703 McLachlan, Philip Diderichsen, and Dorte Haltrup
704 Hansen. 2022. Mending fractured texts. a heuristic
705 procedure for correcting OCR data.

706 Jens Bjerring-Hansen and Matthew Wilkens. 2023.
707 Deep distant reading: The rise of realism in Scandi-
708 navian literature as a case study. *Orbis Litterarum*,
709 78(5):335–352.

710 Jason Dou, Ni Sun, and Xiaojun Zou. 2016. " draw my
711 topics": Find desired topics fast from large scale of
712 corpus. *arXiv preprint arXiv:1602.01428*.

Kenneth Enevoldsen, Lasse Hansen, Dan S. Nielsen,
713 Rasmus A. F. Egebæk, Søren V. Holm, Martin C.
714 Nielsen, Martin Bernstorff, Rasmus Larsen, Pe-
715 ter B. Jørgensen, Malte Højmark-Bertelsen, Pe-
716 ter B. Vahlstrup, Per Møldrup-Dalum, and Kristoffer
717 Nielbo. 2023. [Danish foundation models](#). *Preprint*,
718 *arXiv:2311.07264*. 719

Angela Frattarola. 2018. *Modernist soundscapes: Au-*
720 *ditory technology and the novel*. University Press of
721 Florida. 722

Maarten Grootendorst. 2022. Bertopic: Neural topic
723 modeling with a class-based tf-idf procedure. *arXiv*
724 *preprint arXiv:2203.05794*. 725

Svenja Guhr and Mark Algee-Hewitt. 2024. What’s
726 that scary sound? ambient sound in gothic fiction.
727 *Journal of Computational Literary Studies*, 2(1). 728

Thomas N Haider. 2019. Diachronic topics in new high
729 german poetry. *arXiv preprint arXiv:1909.11189*. 730

Gerhard Heyer, Cathleen Kantner, Andreas Niek-
731 ler, Max Overbeck, and Gregor Wiedemann. 2017.
732 Modeling the dynamics of domain specific ter-
733 minology in diachronic corpora. *arXiv preprint*
734 *arXiv:1707.03255*. 735

Ari Y Kelman. 2010. Rethinking the soundscape: A
736 critical genealogy of a key term in sound studies. *The*
737 *senses and society*, 5(2):212–234. 738

Per E Kummervold, Javier De la Rosa, Freddy Wet-
739 jen, and Svein Arne Brygfjeld. 2021. [Operationaliz-](#)
740 [ing a national digital library: The case for a norwe-](#)
741 [gian transformer model](#). In *Proceedings of the 23rd*
742 *Nordic Conference on Computational Linguistics*
743 *(NoDaLiDa)*, pages 20–29, Reykjavik, Iceland (On-
744 line). Linköping University Electronic Press, Swe-
745 den. 746

Nora Elisabeth Lambrecht. 2017. *The Art of Noise:*
747 *Literature and Disturbance 1900-1940*. Ph.D. thesis,
748 Johns Hopkins University. 749

Jani Marjanen, Elaine Zosa, Simon Hengchen, Lidia
750 Pivovarova, and Mikko Tolonen. 2020. Topic mod-
751 elling discourse dynamics in historical newspapers.
752 *arXiv preprint arXiv:2011.10428*. 753

Jarmila Mildorf. 2019. Can sounds narrate? prosody in
754 sound poetry performance. *CounterText*, 5(3):294–
755 311. 756

Saif Mohammad. 2016. A practical guide to senti-
757 ment annotation: Challenges and solutions. In *Pro-*
758 *ceedings of the 7th workshop on computational ap-*
759 *proaches to subjectivity, sentiment and social media*
760 *analysis*, pages 174–179. 761

Dan Saattrup Nielsen. 2023. Scandeval: A benchmark
762 for scandinavian natural language processing. *arXiv*
763 *preprint arXiv:2304.00906*. 764

David Novak and Matt Sakakeeny. 2015. *Keywords in*
765 *sound*. Duke University Press. 766

- 767 Jakob Ingemann Parby. 2021. Fremskridtets lyd? Lydrevolutionen og håndteringen af støj under københavns industrialisering ca. 1850-1910. *Kulturstudier*, 12(2):41–71.
- 768
- 769
- 770
- 771 John M. Picker. 2003. [3INTRODUCTION: The Tramp of a Fly’s Footstep](#). In *Victorian Soundscapes*. Oxford University Press.
- 772
- 773
- 774 Margaret E Roberts, Brandon M Stewart, Dustin Tingley, Edoardo M Airoidi, et al. 2013. The structural topic model and applied social science. In *Advances in neural information processing systems workshop on topic models: computation, application, and evaluation*, volume 4, pages 1–20. Harrahs and Harveys, Lake Tahoe.
- 775
- 776
- 777
- 778
- 779
- 780
- 781 R Murray Schafer. 1993. *The soundscape: Our sonic environment and the tuning of the world*. Simon and Schuster.
- 782
- 783
- 784 Vésteinn Snæbjarnarson, Annika Simonsen, Goran Glavaš, and Ivan Vulić. 2023. Transfer to a low-resource language via close relatives: The case study on faroese. In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, Tórshavn, Faroe Islands. Linköping University Electronic Press, Sweden.
- 785
- 786
- 787
- 788
- 789
- 790
- 791 Leon Strømberg-Derczynski, Manuel Ciosici, Rebekah Baglini, Morten H. Christiansen, Jacob Aarup Dalsgaard, Riccardo Fusaroli, Peter Juel Henrichsen, Rasmus Hvingelby, Andreas Kirkedal, Alex Speed Kjeldsen, Claus Ladefoged, Finn Årup Nielsen, Jens Madsen, Malte Lau Petersen, Jonathan Hvithamar Rysstrøm, and Daniel Varab. 2021. [The Danish Giga-word corpus](#). In *Proceedings of the 23rd Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 413–421, Reykjavik, Iceland (Online). Linköping University Electronic Press, Sweden.
- 792
- 793
- 794
- 795
- 796
- 797
- 798
- 799
- 800
- 801
- 802 Timothy R Tangherlini and Peter Leonard. 2013. Trawling in the sea of the great unread: Sub-corpus topic modeling and humanities research. *Poetics*, 41(6):725–749.
- 803
- 804
- 805
- 806 Leah Hutchison Toth. 2016. Resonant texts: Sound, noise, and technology in modern literature.
- 807
- 808 Neil Verma. 2019. *Theater of the mind: imagination, aesthetics, and American radio drama*. University of Chicago Press.
- 809
- 810