# CS447 Literature Review: How NLP Powers Music Recommendations: The Role of Neural Summarization and Sentiment Analysis

Felix Koczan,
fkoczan2@illinois.edu

May 5, 2025

### Abstract

Natural-language user feedback can capture the subtle reasons and emotions that underlie listening preferences. Playlist titles, social-media posts or music reviews provide an informative source for listening preferences, yet most recommender systems still rely on numeric ratings or play counts. This literature review asks the question: How can modern NLP techniques enhance music recommendation by extracting insights from such unstructured text? This review examines six papers, three in neural abstractive summarization, Rush et al. (2015), See et al. (2017) and Gehrmann et al. (2018) and three in sentiment and connotation analysis, Buechel et al. (2020), Rashkin et al. (2016) and Rashkin et al. (2017). The works of these authors evolve from attention-based seq2seq models to pointer-generator networks and bottom-up frameworks, each of which improves content fidelity and domain adaptability. The sentiment papers complement this by providing multilingual emotion lexicons and role-based connotation frames that uncover explicit and implicit affect across 91 languages. By combining these approaches, one can generate content-faithful summaries with nuanced, multilingual sentiment cues, improving both item and user profiles. The review highlights ongoing challenges like domain-specific slang and real-time scalability, and proposes future directions such as integrating large language models and incorporating lyrics analysis. Ultimately, it shows that leveraging both user expression and emotion can lead to more empathetic, explainable music recommendation systems.

## 1 Introduction

Music recommendation systems aim to connect listeners with songs and artists they will enjoy. Traditional recommenders often rely on numerical ratings and implicit user behaviors such as plays and likes to model preferences (Al-Ghuribi and Noah, 2021). However, many platforms lack explicit ratings, and users increasingly express opinions through text. This can include comments, reviews, forum posts, and social media posts. These user-generated content sources contain insights that go beyond the plain numerical statistics and tap into the *why* a user likes or dislikes a song. Mentions of melody, lyrics or mood are often coupled with sentiment. By tapping into unstructured textual data, recommendations can uncover more nuanced user preferences that are not evident from ratings alone (Al-Ghuribi and Noah, 2021). This review will look into how Natural Language Processing (NLP) can be utilitzed in music recommender systems. In particular, NLP techniques like neural summarization and sentiment analysis will be analyzed, as these offer powerful tools to automatically extract user insights at scale. The main research question guiding this review is:

*How can NLP enhance the effectiveness of music recommendation systems by leveraging neural summarization and sentiment analysis to extract insights from user-generated content?*

This review addresses the question by examining key NLP advances in abstractive summarization as well as sentiment and connotation analysis, and discussing their relevance to music recommenders. The motivation for this literature review is simple: by summarizing large volumes of user comments and detecting the sentiment and nuanced connotations in user language, a system can better understand the why behind user preferences. For example, if many listeners praise an artists "emotional delivery" or note a song "feels energetic and uplifting," NLP can distill those sentiments and the recommender can use them to suggest music with similar qualities (Dang et al., 2021).

This review is structured as follows. First, a background on abstractive summarization and sentiment analysis, including connotation frameworks, is established. Next, six seminal NLP papers are analyzed: three on neural summarization (Rush et al., 2015); (See et al., 2017); (Gehrmann et al., 2018)) and three on sentiment and connotation analysis (Buechel et al., 2020); (Rashkin et al., 2016); (Rashkin et al., 2017). The contributions of each paper are summarized and connected to the challenge of leveraging user-generated text for music recommendations. Next, the selected works are compared, and their methods and applicability are discussed, with attention to challenges such as domain adaptation and the use of slang. Finally, directions for future work, such as integration into broader recommender systems and handling of domain specific knowledge are outlined. The review concludes with an integrated perspective on how NLP can improve music recommender systems.

## 2   Background

### 2.1   Neural Summarization

Summarization algorithms aim to produce a condensed version of a text that preserves its core meaning. There are two main methods to do so, namely extractive and abstractive summarization. Extractive methods select and concatenate passages from the original text, whereas abstractive methods generate new phrases and potentially rephrase content, similar to how a human would summarize (Rush et al., 2015). Abstractive summarization is more flexible but also more challenging, as it requires natural language generation and an understanding of the source. Early abstractive systems relied on fixed rules, such as deletion-based sentence compression or syntactic transformations. Modern approaches are data-driven, using neural networks to learn how to generate summaries from large corpora of article-summary pairs. In particular, sequence-to-sequence (seq2seq) models with attention mechanisms have become a cornerstone. In seq2seq an encoder network reads the source text and a decoder network generates the summary, while attention provides a learned soft alignment between source and summary (Rush et al., 2015). This approach was inspired by neural machine translation (NMT) and first applied to summarization by Rush et al. (2015). The advantage of neural summarizers is their fluency and trainability on large data, but a crucial drawback is that they face issues like inaccurate reproduction of facts and repetition. To address these issues, the pointer-generator network, which allows copying words from the source, and coverage mechanisms, tracking and avoiding already summarized content, were introduced by See et al. (2017). These also enable summarization of more informal texts. Especially crucial in music recommender systems, where an abstractive summarizer could digest hundreds of user comments about a song into a concise blurb capturing the consensus. Such a summary could give the recommender a quick

understanding of the item's reception and characteristics.

## 2.2   Sentiment Analysis and Connotation

Sentiment analysis is the NLP task of determining the affective attitude expressed in text. This includes, classifying a user review as positive, negative, or neutral. But sentiment analysis can also go beyond these rather simple classifications and analyze emotions, such as specific feelings (e.g. joy, sadness, anger). A fundamental resource for sentiment analysis is the so-called sentiment lexicon (Buechel et al., 2020). A list of words annotated with sentiment polarity or emotion associations. For example, in the lexicon "excellent" is positive, "dismal" is negative. These lexicons are manually curated and thus only exist for a few languages. Buechel et al. (2020) tackled this via a novel method to automatically build multilingual emotion lexicons. The author increased coverage to 91 languages (Buechel et al., 2020). However, user-generated text often is not as explicit in conveying sentiment, as it often contains slang, idioms, or implicit sentiments that are not obvious from single words. This is where connotation analysis comes in. Connotation is commonly understood as the implied sentiment and relationships evoked by language, which go beyond the literal meaning. Rashkin et al. (2016) introduced Connotation Frames, which is a framework to capture the often subtle roles and implications of predicates, mostly verbs. The way connotation frames work is that for a given verb, a connotation frame encodes attributes like the writer's perspective toward the participants, the effect on the entities, and their perceived mental states (Rashkin et al., 2016). For example, saying "X violated Y" implicitly frames X as malicious (negative connotation toward X) and Y as a sympathetic victim. These frames provide a structured way to infer sentiment even if the text does not explicitly use positive or negative adjectives. In user music reviews, connotation frames could detect subtleties: e.g., "This song made my day". The verb made frames the song positively and implies the listener was in a possibly less positive state, which was improved by the song. Such nuanced understanding goes beyond a basic lexicon. Connotation frames were later also extended to multiple languages and applied to social media text, which proves their utility in different contexts (Rashkin et al., 2017).

# 3   Review of Key NLP Papers

## 3.1   Rush et al. (2015): Attention-Based Neural Summarization

In their paper Rush et al. (2015) pioneered the use of neural networks for abstractive summarization (Rush et al., 2015). They proposed an Attention-Based Summarization (ABS) model, which essentially is a seq2seq model with an attention mechanism. The goal was to generate very short headline-style summaries for sentences. At the time of the paper most summarizers were extractive, so the authors wanted to demonstrated that a neural summarization model could learn to paraphrase and compress text. The ABS model uses a local attention technique to generate each word of the summary conditioned on the source sentence. The encoder reads the input into vector representation, and the decoder produces a summary word-by-word, at each step focusing on the most relevant part of the input via attention. The main benefit of Rush et al. (2015) approach was that it was structurally simpler than many of the prior abstractive systems, yet it was effective and trainable on large data. The authors trained their model on a corpus of news articles (around 4 million examples) and evaluated on the DUC-2004 shared task, which is a standard benchmark for summarization. The result was significant performance gains over previous baselines. The ABS model outperformed

classical approaches and even some systems that relied on heavy manual engineering (Rush et al., 2015).

Rush et al. (2015) work is important in that it proved that neural networks with attention can perform abstractive summarization. For the domain of music recommenders, this means that it is feasible to automatically generate concise summaries of user comments about music. The ABS model could, for instance, produce a one-line summary of a long user review highlighting the key opinion (e.g. "I found the album all over the place, really loved two songs, but the rest felt like filler material" into a headline like "Album has a few great tracks but lacks consistency"). By capturing the core idea, a recommender system can figure this user's nuanced view. One limitation of Rush et al. (2015) original model is that it was designed for sentence-level summaries such as headlines, not multi-sentence inputs. In a music recommender scenario, we often want to summarize a multitude of comments or a multi-sentence review. Nonetheless, the techniques introduced laid the groundwork for later models that handle longer inputs. Rush et al. (2015) model also incorporated a simple mechanism to include certain keywords, which already hinted at later developments like the pointer mechanism. In summary, this paper's contribution is the proof-of-concept that neural abstractive summarization is viable and efficient (Rush et al., 2015).

### 3.2   See et al. (2017): Pointer-Generator Networks for Summarization

See et al. (2017) extended neural summarization with two key innovations (See et al., 2017). Their model, which is often called the Pointer-Generator Network, augments the seq2seq attentional model with a pointer mechanism and a coverage mechanism. The pointer mechanism allows the decoder to either generate words from its vocabulary or copy words from the source text via pointing. This hybrid approach is important for dealing with out-of-vocabulary words and proper nouns. Especially in music reviews, the model could just copy the name of an artist or song title from the user's comment rather than producing an unknown token. This helps ensure the accurate reproduction of information that should remain consistent over time (i.e., factual content), while still allowing the model to generate new words when needed. The second key innovation in their paper is the coverage mechanism, which tracks attention over the source text to discourage the decoder from repeatedly focusing on the same content. This addresses a common issue in neural summarization, the generation of repeated phrases or redundant sentences (See et al., 2017).

In their paper, See et al. (2017) applied their pointer-generator with coverage model to the CNN/Daily Mail news article dataset and achieved a new state-of-the-art performance. It outperformed the previous best abstractive system by > 2 ROUGE points. The summaries produced by their model were less prone to missing important facts and had far fewer repeated phrases (See et al., 2017). These can be attributed to the advances by the pointer and coverage mechanism introduced in the paper.

See et al. (2017) work is directly relevant to summarizing user-generated content for recommendations. User reviews and comments about music often include unusual words. This could be slang, misspellings, or specific references to the music world, such as artist names or genre labels. A pointer-generator model could handle these well by just copying from the source text, ensuring important terms are carried into the summary (See et al., 2017). For example, if many users mention "lo-fi aesthetic" or a particular subgenre, a simple vanilla neural summarizer might not have those in vocabulary, but a pointer mechanism can copy them so the summary retains those key phrases. Coverage is equally valuable in music recommendations. User comments can be verbose or go on a tangent. The coverage mechanism will systematically cover each point without looping back to already covered content. The twofold contribution of See et al. (2017) translates to more reliable

extraction of what listeners are saying (See et al., 2017). One critique of pointer-generators is that they tend to be "poorly abstractive," leaning heavily on copying. In other words, they sometimes just string together pieces of the source text rather than truly rephrasing (Scialom et al., 2020). This is a trade-off. For a music recommender, preserving original wording is actually a plus, but it might result in less generalization. Nonetheless, See et al. (2017) model remains a robust choice for domains like user reviews, and its design ideas (copying and coverage) have influenced nearly all subsequent summarization systems.

### 3.3 Gehrmann et al. (2018): Bottom-Up Abstractive Summarization

While seq2seq models generate summaries in a single pass, Gehrmann et al. (2018) proposed a two-step bottom-up approach to abstractive summarization (Gehrmann et al., 2018). The problem they targeted was the content selection issue. Neural models sometimes include unnecessary details or even copy full sentences from the source, which limits usefulness. Gehrmann et al. (2018) solution was to introduce a content selector that operates before the neural generator. This content selector learns to identify phrases or segments in the source document that are likely to be part of the summary. By over-determining the important phrases in the source text, the model creates a guided roadmap for the generator. In the second step, an abstractive summarization model, much similar to prior seq2seq models, generates the summary but is constrained via this predefined bottom-up attention to focus on the selected phrases (Gehrmann et al., 2018). Essentially, it can only attend to and include those high-importance portions of the input, ensuring that less relevant content is ignored.

The authors two-step bottom-up process was shown to improve the relevance of summaries. They reported that their method produced more concise summaries, with fewer unnecessary bits, while maintaining fluency. In terms of results, it outperformed end-to-end models on standard benchmarks, achieving higher ROUGE scores on both the CNN/DailyMail and New York Times datasets. An important advantage of the content selector is its efficiency. It can be trained with as little as 1,000 sentences, making it very data-efficient even for new domains (Gehrmann et al., 2018). This means one can adapt a summarizer to a different genre or domain by just labeling a small set of summaries for that domain to train the selector on, rather than retraining a whole seq2seq model from scratch.

The bottom-up approach described by (Gehrmann et al., 2018) is highly relevant for leveraging user-generated content in music recommendations because it directly addresses the two critical challenges of domain adaptation and content focus. User discussions about music might include tangential content, such as off-topic conversations or unrelated anecdotes that a generic summarizer could mistakenly include. A content selector trained on a bit of in-domain data could relatively easily learn to pick out just the musically relevant phrases and ignore irrelevant chatter. Gehrmann et al. (2018) found that only about 1k training examples suffice for the selector is encouraging, since there may not be a massive labeled dataset of music review summaries readily available. The real important aspect is the transferability (Gehrmann et al., 2018). One could take a summarization model pre-trained on news and quickly tune it to summarize music-related texts by tuning a content selector. Gehrmann et al. (2018) work contributes a means to make summarization more accurate and domain-aware, which is crucial in the noisy, varied world of user content.

### 3.4 Buechel et al. (2020): Learning Multilingual Emotion Lexicons

Buechel et al. (2020) focused on a fundamental resource gap in sentiment analysis. The lack of

emotion lexicons for most languages. Their work, "Learning and Evaluating Emotion Lexicons for 91 Languages," (Buechel et al., 2020) introduced a methodology to automatically build large lexicons of emotion-linked terms in many languages using minimal resources. The approach described in their paper requires only an existing lexicon in one source language, which in practice is most often English, a bilingual translation model, and a word embedding model for the target language. So to put simply, they project emotion knowledge from English to other languages by translating the words and then refining using distributional similarity. The authors did this for 91 languages and thus generated lexicons that each contain over 100,000 entries, covering eight emotional variables. The emotional variables included cover a broad spectrum. For example, they mention eight variables encompassing things like basic emotions such as joy or anger and dimensions like valence, arousal and dominance (Buechel et al., 2020).

Using human-labeled emotion datasets in 12 diverse languages, they found the quality of the generated lexicons to be on par with state-of-the-art monolingual lexicon creation methods. In some cases, the auto-generated lexicons even exceeded human annotator reliability on certain emotions, showing that the method can capture nuances (Buechel et al., 2020).

In a music recommendation scenario, especially one with a multilingual user base, this work provides the building blocks to perform sentiment and emotion analysis on user comments in many languages, beyond English. For example, one might consider a song that's popular in Latin America, thus many reviews might be in Spanish or Portuguese. With Buechel et al. (2020) lexicons, one could identify emotion-laden words in those reviews and quantify the sentiment. The lexicons cover not just positive/negative polarity but a range of emotions, which means the system can detect the mood a song elicits. This is crucial to enhance recommendations by matching songs to users' current mood preferences. Buechel et al. (2020) effectively removed a barrier to multilingual sentiment analysis. One critique is that lexicon-based analysis alone can miss context such as sarcasm or negation. In summary, Buechel et al. (2020) contribute a multilingual foundation for sentiment detection, which is essential for inclusive recommendation systems.

### 3.5 Rashkin et al. (2016 & 2017): Connotation Frames for Implied Sentiment

Rashkin et al. (2016) work on Connotation Frames and its extension, Rashkin et al. (2017) focus on implied sentiments and relationships. In their 2016 ACL paper, Rashkin, Singh, and Choi introduced connotation frames as a formalism to represent the latent affective content of verbs. Unlike traditional sentiment analysis which might label a whole sentence positive or negative, connotation frames break a sentence down into more nuance. They assign perceived attitudes to the participants in an event and to the speaker. Specifically for a transitive verb, the connotation frame might include attributes such as the writer's perspective toward the subject and object, the perspectives of the subject and object toward each other, the implied value of each entity, and the effect or mental state experienced by each (Rashkin et al., 2016). For example, in the case of the verb 'violate', the connotation frame suggests that the writer views the subject (X) negatively, as an antagonist and the object (Y) positively, as a victim. It also implies that Y is something valuable (since violating is framed as harmful), and that Y experiences a negative outcome, such as harm or distress. More concretely, in a sentence like 'The company violated the user's privacy,' the company (X) is viewed unfavorably, the user (Y) is seen as a victim, and privacy is implied to be something valuable and unjustly harmed. Rashkin et al. (2016) constructed a lexicon of connotation frames for around 950 common verbs and then built predictive models to infer these frames from distributional semantics. They showed that it is possible to automatically predict the connotative relations of verbs, and they further demonstrated an application in analyzing bias in news media in their paper (Rashkin et al.,

2016).

In 2017, Rashkin and collaborators extended their idea on connotation frames (Rashkin et al., 2017). In their 2017 paper, they created Multilingual Connotation Frames and applied them to a social media context. They expanded the lexicon to 10 languages and used the frames to analyze targeted sentiment on Twitter around public events. They used 1.2 million extracted frames from Twitter to forecast shifts in public sentiment towards entities, managing to predict changes up to half a week in advance. This showed the approach's power in capturing dynamic sentiment signals in noisy, informal text across languages (Rashkin et al., 2017).

Connotation frames are a more granular and context-rich sentiment analysis tool that could greatly benefit a music recommender's understanding of user comments. Much of user feedback is not plain and obvious such as "I like this" or "I hate this". Often sentiment is conveyed through storytelling or metaphors which essentially are connotations. For instance, a user might say "This song hits hard" (positive connotation toward the song, implying it positively affects the listener), or "The album bombed for me" (using "bomb" in a negative sense, implying a negative sentiment). A standard sentiment model might misclassify such cases without additional knowledge. Connotation frames, provide a way to infer the implied sentiment roles. By incorporating a connotation lexicon or model, the system can catch these subtleties better than by word-level sentiment alone.

The multilingual aspect from 2017 is particularly relevant if users around the world use different idioms to convey sentiment. Rashkin et al. (2017) multilingual frames show it's possible to extend these nuanced analyses beyond English. Also, the fact they applied it to targeted sentiment analysis on social media suggests a use case for recommendation: monitoring how sentiment toward a new song evolves over time across different listener groups on social media. A music recommender could use such signals (e.g., connotation frames indicating growing positive sentiment in a certain region or demographic) to adjust recommendations in near-real-time.

In summary, connotation frames add a layer of depth to sentiment analysis, capturing the attitudes and connotative implications behind what people write. Overall, Rashkin et al.'s contributions enable sentiment analysis to go beyond surface-level word sentiment to contextual, role-based sentiment inference.

## 4    Comparative Analysis of Approaches

The six papers reviewed span two interconnected areas: abstractive/neural summarization (Rush et al. (2015); See et al. (2017); Gehrmann et al. (2018) and sentiment/emotion analysis Buechel et al. (2020); Rashkin et al. (2016); Rashkin et al. (2017). The reviews of each paper reveal a complementary relationship: the summarization works focus on extracting and compressing information, while the sentiment-focused papers interpret the meaning behind that information. Together, they address both the *what* and the *why* of user-generated text.

### 4.1    Summarization Models - Evolution and Content Fidelity

Rush et al. (2015) introduced the baseline neural approach (ABS) that proved the feasibility of abstractive summarization (Rush et al., 2015). This was a single-step encoder-decoder with attention, effective for short inputs. See et al. (2017) built on this by tackling practical issues like out-of-vocabulary words and repetitive outputs through the pointer-generator and coverage mechanisms (See et al., 2017). This made summaries more accurate to the source (copying names, facts) and cleaner (no loops). Gehrmann et al. (2018) then added a content selection step to explicitly ensure

the summarizer focuses on the right parts. This evolution is important for user-generated content: early neural models might generate plausible-sounding but irrelevant summaries if users wander off-topic (Gehrmann et al., 2018). Later models enforce staying true to salient content. One common thread is that all three models were demonstrated on news articles (with Rush and See using CNN/DailyMail or DUC, and Gehrmann on CNN/DM and NYT). News is well-structured. User text is messier. Thus, while these models provide a strong starting point, they likely need adaptation for the domain of music recommendations. Gehrmann et al. (2018) bottom-up approach explicitly addresses adaptation by using small in-domain data for the selector, which is a notable benefit (Gehrmann et al., 2018). See et al. (2017) pointer mechanism is another domain-friendly feature, possibly to be used for the slang and proper nouns in music discussion. In summary, all three summarization papers aim to maximize relevant content and coherence in generated summaries, with varying strategies: Rush et al. (2015) rely on the neural model's implicit learning, See et al. (2017) add internal fixes (copy & coverage) to handle accuracy, and Gehrmann et al. (2018) add an external content filter for focus.

## 4.2   Sentiment Analysis Advances – Lexicons vs. Connotation Frames

The two sentiment-focused papers represent two different approaches. One lexicon-based and expansive (Buechel et al., 2020), the other frame-based and contextual (Rashkin et al., 2016). Buechel et al. (2020) is about breadth. Scaling emotion lexicons to many languages automatically. It follows a more traditional view that words carry sentiments and emotions that can be catalogued. Rashkin et al. (2016); Rashkin et al. (2017) is more about depth. Understanding that sentiment can be conveyed through the interplay of words and may not be so explicit. Lexicons (even big ones) might label words like "great" or "terrible," but they will not directly tell you that "X broke Y's heart" implies a negative sentiment toward X. Connotation frames fill that gap by providing a framework for implied sentiment roles. When applicable, frames give a much fuller picture. Who is viewed positively, who negatively, and what the nuanced effect is (Rashkin et al., 2016). For a music recommender, this means lexicons can quickly gauge overall sentiment (e.g., 80% of words in this review are positive, so the review is positive), whereas connotation frames could catch something like "The song leaves the audience in tears" (which by frames might imply the audience is sad but perhaps in a good way if context is emotional performance). Another distinction is evaluation domain. Buechel et al. (2020) validated on emotion datasets and lexicon quality intrinsically, while Rashkin et al. (2016) applied frames to news bias and Twitter sentiment forecasting (Buechel et al., 2020); (Rashkin et al., 2016). The latter shows the power in real-world scenarios, indicating that frames are effective in analytical and predictive tasks. This suggests that for music, connotation frames could help predict a song's reception trajectory.

In summary, the common goal of all six papers is to extract meaningful information from text. Summarization papers extract informative nuggets and gist, sentiment papers extract affective and subjective signals. Both types are crucial for answering the proposed research question. They also complement each other. Summarization can condense a long discussion into key points, and sentiment analysis can then be applied to those points to gauge positive/negative valence. Or vice versa, sentiment analysis might identify the most emotionally charged comments, which a summarizer could focus on.

# 5 Discussion

The reviewed approaches are highly synergistic for the posed problem. In practice, a music recommendation system could and would employ summarization to distill the crowd's opinions on each song/album and sentiment analysis to quantify those opinions. For example, one could imagine a new album release with thousands of comments. An abstractive summarizer could produce a summary such as "Listeners praise the album's playful lyrics, though some found the middle tracks repetitive." This summary provides a quick overview of opinions. Then, using sentiment analysis, the system can tag this summary (or the underlying comments) with sentiments. Playful lyrics = positive, repetitive = negative. From connotation frames, the verb "praise" indicates a positive writer perspective towards the album, while "found X repetitive" suggests a negative evaluation of X. By combining these, the recommender knows the album is generally liked for specific reasons, with a minor critique. This can improve recommendation in three ways.

1. User Preference Matching. If a user frequently expresses liking for "playful lyrics" in other music, the system can match them to this album which is summarized as having playful lyrics.

2. Item Profiling. The summary and sentiment tags form a profile of the item, beyond metadata like genre. This helps in cold-start for the item. Even if the system has little listening data, the text-derived profile guides who might like it (e.g. people who like songs with playful lyrics).

3. Explanations. Showing the summary to users as an explanation ("People say XYZ") can guide new listeners and increases transparency of the recommendation.

## 5.1 Strengths and Limitations

A strength of neural summarization models is their fluency and adaptability (See et al., 2017); (Gehrmann et al., 2018). They can generate human-like summaries tailored to what a user might care about (with bottom-up guidance). However, they are not perfect. One critique is that seq2seq summarizers can sometimes produce factually inconsistent or generic statements if the input is out-of-distribution. User comments might include slang or unconventional grammar that confuses a model trained on news, which most likely would not include such words. The pointer-generator helps here by copying unknown words, but it does not guarantee understanding of those words. There is also a risk of loss of nuance. Sarcasm for instance in a user comment could be lost in summary. Also, early neural summarizers tended to optimized for very short summaries. Music reviews most likely need longer summaries to capture multiple aspects. Bottom-up can handle moderate lengths (they were used for multi-sentence summaries), but the content selector will need careful tuning to not omit too much. The coverage mechanism from See et al. (2017) directly addresses redundancy, but it does not fully solve the problem if the model learns to be extractive. Some later research noted pointer-generators still copy large chunks of text (Scialom et al., 2020). In our use case, copying is double-edged, because copying an important phrase is great (we want "heart-wrenching" if that is what everyone says about a song), but copying whole sentences from a user comment could introduce irrelevant personal remarks.

On the sentiment side, lexicon-based analysis Buechel et al. (2020) is lightweight and interpretable (Buechel et al., 2020). One can easily see which words contributed to a sentiment score. But one limitation is handling of context and composition. If a user says "not bad" or uses irony or sarcasm, a lexicon might be misled ("bad" is negative, but "not bad" means okay or good). Simi-

larly, lyrics or artist names could be misidentified as sentiment (e.g. for the band "The Killers" the word "Killer" might receive a negative sentiment).

Connotation frames offer a more nuanced context handling, but they are inherently limited to certain constructs. They need verbs with clear agent/theme. A lot of music commentary might be adjective-based ("this track is fire!") or noun-based ("an instant classic"), where connotation frames do not directly apply. Additionally, deploying connotation frame models can be computationally heavier. One must parse sentences, identify verbs and arguments, then apply the frame predictions (Rashkin et al., 2016). It is more involved than a simple lexicon lookup, but the payoff is potentially higher precision in understanding complex opinions.

## 5.2 Applicability

All six papers reported strong results in their domains, often state-of-the-art. For recommendations, however, the ultimate metric is not ROUGE or lexicon accuracy but whether using these NLP techniques improves recommendation quality, which might be measured in accuracy and ultimately user satisfaction. There have been studies demonstrating that incorporating sentiment from reviews improves recommender performance (Dang et al., 2021); (Barrière and Kembellec, 2018). Those successes can be interpreted as validation that the general idea works. Understanding what users say and how they feel about items leads to better modeling of user preferences. While the papers reviewed in this literature provide advanced methods for understanding user text, deploying these models in a live system requires substantial adaptation and engineering effort. Summarization models need to be efficient to run on potentially millions of items and in near real time for new reviews. Pointer-generator networks and content selectors can be resource-intensive, because essentially they are neural models with many parameters. One might consider using them offline to periodically generate summaries for each item, which are then stored. Similarly, sentiment lexicons and connotation models could be used offline to tag items and users with sentiment signals that feed into the recommendation algorithm.

In conclusion of this discussion, the NLP methods reviewed offer a toolkit for music recommender systems. Sequence-to-sequence models (with pointing and coverage) ensure we capture what listeners talk about, and lexicons/frames ensure we capture how they talk about it (positive, negative, as well as implied attitudes). Each has limitations, but combined, they significantly push the boundaries of what a recommender system can understand from text.

# 6 Future Work

As can be seen from this review, while promising, applying neural summarization and sentiment analysis to music recommendation raises several challenges and open research questions.

## 6.1 Domain Adaptation

Summarization models trained on news need adaptation to the music domain's content and style. As Gehrmann et al. (2018) showed, even only 1,000 in-domain examples can help (Gehrmann et al., 2018). Future work could focus on transfer learning for summarization. For example, fine-tuning a pre-trained summarizer on a small set of music reviews paired with expert-written summaries. These might be taken from music blogs or summary-style reviews. Techniques like data augmentation could generate pseudo-examples of music talk to expand training data. Domain adaptation also entails adjusting to the shorter, often informal structure of user comments compared to long

news articles. Research into summarizing tweets and forum posts can be leveraged, possibly by combining summarization with techniques for noisy text normalization.

## 6.2   Slang and Figurative Language

Music fans often use creative and highly context-dependent language, for example, phrases like "this track is lit," "an earworm," "fire beats," or "kills it on stage." Many of these expressions are idiomatic and cannot be interpreted or translated literally. A pointer-generator can copy "lit" into a summary, but understanding it requires either a dedicated "slang"-lexicon entry or context. Future sentiment models might integrate Urban Dictionary style resources or use context from social media to learn such slang sentiment. Metaphors (such as "an earworm") are another challenge. This may require extending connotation frames or using figurative language processing techniques. An interesting direction might be to develop an irony/sarcasm detector to enhance sentiment analysis, so that a comment like "Great, another song to put me to sleep" is not wrongly tagged as positive about inducing sleep, but correctly seen as sarcasm.

## 6.3   Integration into Recommendation Algorithms

Another area for future research is developing models that smoothly integrate text-derived features with collaborative filtering. There is room for end-to-end approaches where, say, a transformer could read all reviews and listening history jointly. However, that might be impractical at scale. Instead, a pipeline where NLP provides intermediate outputs (summaries, sentiment scores, aspect ratings) which are then used as inputs to a recommendation model seems more practical and modular. How to weight these textual insights against behavioral data remains an open question though. Too high and one might recommend an item that people talk about a lot but actually don't listen to as much, too low and one might miss the value of user generated text. Adaptive methods that learn when text is most useful could be explored.

In addressing these challenges, future research will likely draw on advancements in NLP technologies, such as transformer-based summarizers, large pre-trained language models for sentiment analysis, and domain adaptation techniques. The unique context of music, where both content (e.g., audio and lyrics) and culture (e.g., fan language) play crucial roles, presents interesting opportunities for interdisciplinary innovation at the intersection of natural language processing and computational social science.

# 7   Conclusion

Natural Language Processing provides powerful tools to make music recommendation systems more relevant and user-aware by uncovering what is hidden in unstructured text. This literature review examined how neural summarization techniques can condense user-generated content into succinct, informative representations, and how sentiment and connotation analysis techniques can gauge the affective nuances of user opinions. The six analyzed papers illustrate key advancements over time. From the first neural summarizer by Rush et al. (2015) to the more accurate pointer-generator of See et al. (2017) and the content-guided framework of Gehrmann et al. (2018), one can see that neural summarization has become more precise, controllable and produces less irrelevant output. In parallel, Buechel et al. (2020) expanded the capability to analyze sentiments across languages.

Rashkin et al. (2016); Rashkin et al. (2017) provide tools to read between the lines of user language, in order to capture any sentiment or bias that is implicit.

When applied to music recommendation, these NLP techniques enable a system to not only know that a user played Song X 20 times, but also to understand why the user liked Song X and how they felt about it. This more in depth understanding can enhance the recommendation algorithm's decisions, making them more accurately match thematic or emotional preferences. Moreover, by leveraging multilingual lexicons and connotation frames, the system can become globally aware and extract insights from diverse communities of fans around the world.

Naturally, there remain challenges in bringing these research ideas to life. Adapting models to informal language and scaling to real-time. But the trajectory is clear. As NLP continues to improve, especially with large pre-trained language models now capable of summarization and sentiment tasks with minimal supervision, the integration of textual intelligence into recommender systems will become ever more feasible. Ultimately, the goal is to enhance personalization by listening not only to the music itself, but also to what listeners say about it, enabling recommendation systems to capture a deeper understanding of user preferences. They can recommend songs that resonate with listeners not only because "people who liked A also liked B," but because they understand that "people who *feel* the same way about A, also tend to love B." This kind of insight will help make music recommendations more empathetic and truly user-centric.

# References

Sumaia M. Al-Ghuribi and Shahrul Azman Mohd Noah. 2021. A comprehensive overview of recommender system and sentiment analysis. *arXiv preprint arXiv:2109.08794*.

Valentin Barrière and Gérald Kembellec. 2018. A short review of sentiment-based recommender systems. In *Proceedings of the Digital Tools & Uses Congress*, pages 1–4.

Sven Buechel, Susanna Rücker, and Udo Hahn. 2020. Learning and evaluating emotion lexicons for 91 languages. In *Proceedings of ACL*, pages 1202–1217.

Cach Ngoc Dang, María Natalia Moreno-García, and Fernando De la Prieta. 2021. An approach to integrating sentiment analysis into recommender systems. *Sensors*, 21(16):5666.

Sebastian Gehrmann, Yuntian Deng, and Alexander M. Rush. 2018. Bottom-up abstractive summarization. In *Proceedings of EMNLP*, pages 4098–4109.

Hannah Rashkin, Eric Bell, Yejin Choi, and Svitlana Volkova. 2017. Multilingual connotation frames: A case study on social media for targeted sentiment analysis and forecast. In *Proceedings of ACL (Short Papers)*, pages 673–678.

Hannah Rashkin, Maarten Singh, and Yejin Choi. 2016. Connotation frames: A data-driven investigation. In *Proceedings of ACL*, pages 311–321.

Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. In *Proceedings of EMNLP*, pages 379–389.

Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2020. MLSUM: The multilingual summarization corpus. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8051–8067, Online. Association for Computational Linguistics.

Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of ACL*, pages 1073–1083.