

HapticLLaMA : A Multimodal Sensory Language Model for Haptic Captioning

Guimin Hu¹, Daniel Hershcovich¹, Hasti Seifi²

¹University of Copenhagen

²Arizona State University

rice.hu.x@gmail.com dh@di.ku.dk hasti.seifi@asu.edu

Abstract

Haptic captioning is the task of generating natural language descriptions from haptic signals, such as vibrations, for use in virtual reality, accessibility, and rehabilitation applications. While previous multimodal research has focused primarily on vision and audio, haptic signals for the sense of touch remain under-explored. To address this gap, we formalize the haptic captioning task and propose HapticLLaMA, a multimodal sensory language model that interprets vibration signals into descriptions in a given sensory, emotional, or associative category. We investigate two types of haptic tokenizers, a frequency-based tokenizer and an EnCodec-based tokenizer, that convert haptic signals into sequences of discrete units, enabling their integration with the LLaMA model. HapticLLaMA is trained in two stages: (1) supervised fine-tuning using the LLaMA architecture with LoRA-based adaptation, and (2) fine-tuning via reinforcement learning from human feedback (RLHF). We assess HapticLLaMA’s captioning performance using both automated n-gram metrics and human evaluation. HapticLLaMA demonstrates strong capability in interpreting haptic vibration signals, achieving a METEOR score of 59.98 and a BLEU-4 score of 32.06 respectively. Additionally, over 61% of the generated captions received human ratings above 3.5 on a 7-point scale, with RLHF yielding a 10% improvement in the overall rating distribution, indicating stronger alignment with human haptic perception. These findings highlight the potential of large language models to process and adapt to sensory data.

1 Introduction

Humans perceive their environment through five primary senses: vision, hearing, touch (haptics), taste, and smell. Integrating these sensory modalities can enhance AI systems’ ability to interpret human perception and behavior by providing a

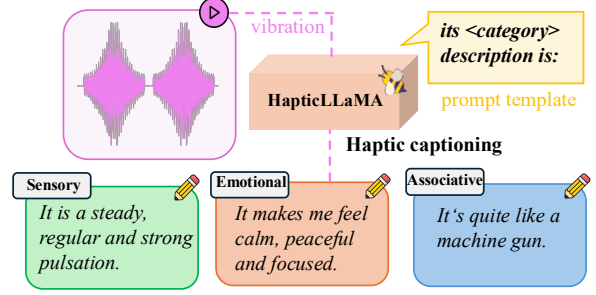


Figure 1: HapticLLaMA can generate sensory, emotional, and associative captions for an input vibration haptic signal with a tailored prompt template.

richer and more nuanced understanding of user context and sensory experiences. Haptic signals convey diverse information through tactile sensations, including physical attributes such as surface texture (Culbertson et al., 2014), emotional states like urgency or pleasantness (Yoo et al., 2015), and recognizable real-world cues such as heartbeats or buzzing of a bee (Seifi et al., 2015; Seifi and MacLean, 2017). These tactile interactions have broad applications in areas including user interactions in virtual reality (VR), physical rehabilitation, blind user navigation, and gaming (Choi et al., 2021; Seim et al., 2022; Katzschmann et al., 2018; Yun et al., 2023a; MacLean et al., 2017).

Multimodal Large Language Models (MLLMs) have significantly advanced captioning across various modalities such as images (Vinyals et al., 2015; Mokady et al., 2021), video (Iashin and Rahtu, 2020; Wu et al., 2023), and audio (Zhang et al., 2022; Liu et al., 2022). *Haptic captioning* involves generating textual descriptions (captions) of haptic signals, capturing sensory, emotional, and associative qualities (see Figure 1). This process enhances the interpretability of physical signals by describing human perception of vibration patterns, similar to how vision-language models capture human perception of images. Additionally, haptic captioning

offers a novel opportunity to develop datasets and benchmarks that probe the limits of AI capabilities, challenging models to accurately interpret and articulate complex physical sensations through natural language.

Unlike image and audio captioning, haptic captioning remains an underexplored domain. As a nascent field, haptic-language understanding presents unique challenges but also significant potential for advancing multimodal AI. The task involves two key challenges: (1) the lack of established tokenization methods for representing haptic signals in a format suitable for language models; and (2) the absence of sensory multimodal models capable of processing haptic vibrations. To date, little work has investigated the ability of large language models (LLMs) to interpret haptic signals.

To support the haptic captioning task, we introduce **HapticLLaMA**, the first multimodal haptic language model trained with two distinct haptic tokenizers that integrates haptic vibration signals and textual descriptions within a single framework. Raw haptic signals are continuous time-series data, incompatible with the discrete token-based input required by LLMs. To address this, we investigate two types of haptic tokenizers: a frequency-based tokenizer and an EnCodec-based tokenizer (Défossez et al., 2023), both designed to convert vibrations into interpretable token sequences suitable for LLMs. HapticLLaMA is trained in two stages: (1) supervised fine-tuning using a LoRA-adapted (Hu et al., 2022) LLaMA architecture (Touvron et al., 2023), and (2) reinforcement learning from human feedback (RLHF; Ouyang et al., 2022) applied using generated haptic captions rated by humans.

Specifically, we contribute: (1) HapticLLaMA, the first haptic language model capable of generating sensory, emotional, and associative captions for vibration signals, (2) two haptic tokenizers that convert raw vibration inputs into sequences of discrete tokens, (3) the VibRate dataset, containing 16,896 user-rated <vibration, caption, rating> samples, used to further incorporate human perception into the model via RLHF, and (4) extensive experiments using automated n-gram metrics and human evaluation demonstrating that HapticLLaMA exhibits strong perceptual capabilities in describing vibration signals, achieving a METEOR score of 59.98 and a BLEU-4 score of 32.06, and receiving human ratings above 3.5 (on a 7-point scale) for over 61% of the generated captions.

2 Related Work

2.1 Haptic Modality and Touch Datasets

Early research on haptic language understanding (Obrist et al., 2013; Seifi et al., 2015; Seifi and MacLean, 2017; Dalsgaard et al., 2022) emerged within the human-computer interaction (HCI) domain, primarily through qualitative studies. While these works underscored the importance of understanding how users describe haptic experiences, they were limited in scale, typically focusing on fewer than 20 signals, and relied on manual analysis methods. Recent work proposed large-scale touch-language datasets for robotic sensing. Datasets such as TVL (Fu et al., 2024), TLV (Cheng et al., 2024b), and Touch100k (Cheng et al., 2024a) use tactile images captured by deformable RGB-like sensors (Yuan et al., 2017). Others collected camera and touch sensor data from human interactions with objects (Balasubramanian et al., 2024; Yang et al., 2022). These datasets convey object properties like shape, size, and texture. While useful for training robots to perceive physical objects, these datasets lack programmable haptic feedback like vibrations from phones and VR controllers, which are common in user-facing applications. Recently, Hu et al. (2024) developed a pipeline for mapping emotional tags to haptic features, but it was only tested on 32 signals with 12 descriptions each. In contrast, HapticLLaMA is the first sensory LLM to generate natural language captions for vibrations, highlighting the broad range of human experiences related to sensory, emotional, and associative aspects of vibration perception.

2.2 Multimodal Captioning

Multimodal models for image (Vinyals et al., 2015; Mokady et al., 2021), video (Iashin and Rahtu, 2020; Wu et al., 2023), and audio captioning (Zhang et al., 2022; Liu et al., 2022) have rapidly advanced in the last decade. Recent advances in open-source LLMs such as LLaMA (Touvron et al., 2023) have accelerated the development of multimodal models. Parameter-Efficient Fine-Tuning (PEFT) (Xu et al., 2023; Xin et al., 2024) adapts pretrained models by updating a small subset of parameters or lightweight modules, significantly reducing computational and storage costs. PEFT has been widely adopted in recent studies (Cheng et al., 2024c; Gema et al., 2023). In the field of haptic-language understanding, Hu et al. (2025) recently introduced HapticCap, a dataset

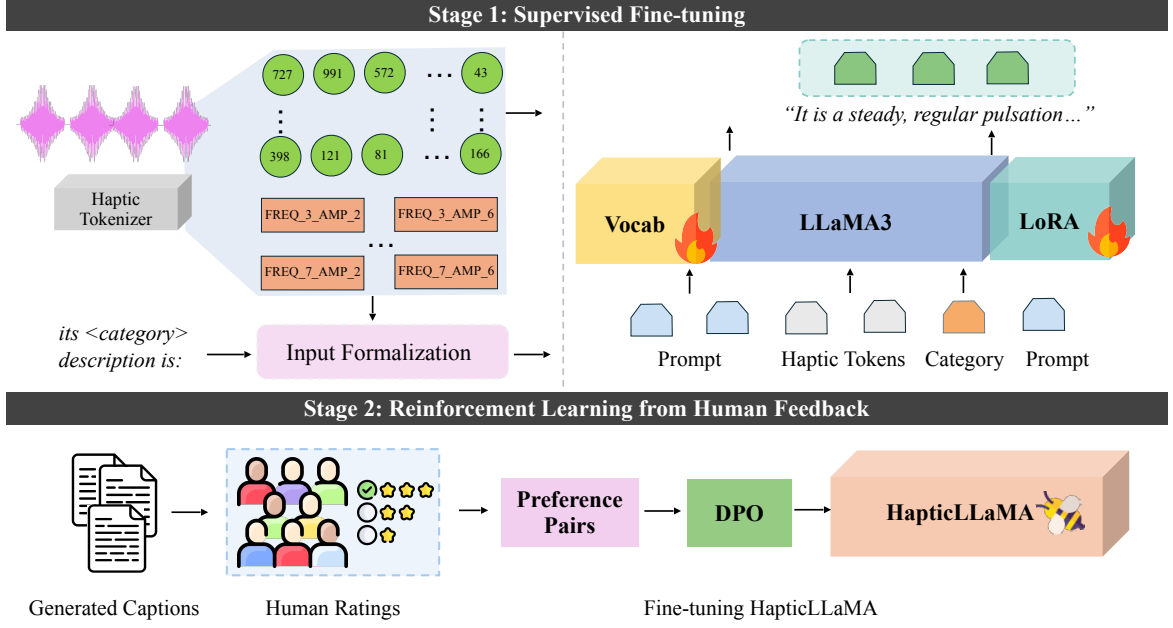


Figure 2: Overview of the two-stage process for constructing HapticLLaMA: (1) supervised fine-tuning, and (2) reinforcement learning from human feedback (RLHF).

of 92,070 vibration samples paired with sensory, emotional, and associative descriptions, laying the foundation for sensory language models in haptics. In our work, we use HapticCap and PEFT to train the first stage of HapticLLaMA, and in the second stage, we incorporate reinforcement learning (Mnih et al., 2015) guided by human feedback.

3 HapticLLaMA

3.1 Task Definition

Similar to image (Vinyals et al., 2015; Mokady et al., 2021) and audio captioning (Zhang et al., 2022; Liu et al., 2022), haptic captioning involves generating textual descriptions (i.e., captions) from haptic signals. Given a vibration signal S and a target category $c \in \{\text{sensory, emotional, associative}\}$, where sensory refers to physical attributes (e.g., intensity of tapping), emotional denotes affective impressions (e.g., the mood of a scene), and associative indicates real-world familiar experiences (e.g., buzzing of a bee, a heartbeat), the goal is to generate a caption corresponding to the specified category of haptic experience.

3.2 Overall Architecture

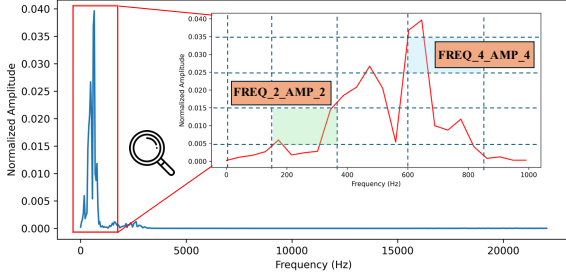
As shown in Figure 2, HapticLLaMA is built on the LLaMA architecture and consists of a haptic tokenizer, a LLaMA model enhanced with LoRA, and a human feedback module trained using rein-

forcement learning. We process the haptic signals offline, converting them into sequences of discrete tokens using two methods: a frequency-based tokenizer using spectral frequency information and an EnCodec-based pretrained neural audio codec (Défossez et al., 2023). The haptic tokens and target category are formatted into a multimodal prompt and input to a LLaMA model fine-tuned with LoRA. Following supervised fine-tuning, human ratings of the generated captions are collected to further refine the model via Direct Preference Optimization (DPO; Rafailov et al., 2023).

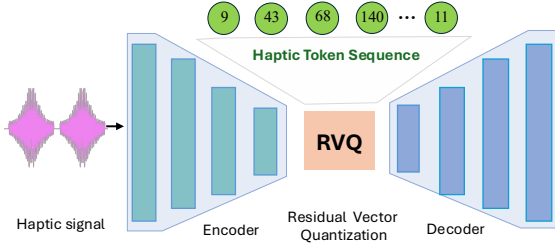
3.3 Input Formalization

3.3.1 Haptic Tokenizer

Frequency-based Tokenizer: This tokenizer (Figure 3(a)) is motivated by the importance of spectral frequency information in characterizing vibrations (Jones and Sarter, 2008; Bensmaïa et al., 2005). Unlike the time domain, where signals are represented by amplitude values over time, the frequency domain represents a signal as a sum of sinusoids at varying frequencies. The frequency-based tokenizer begins by converting the time-domain signal into the frequency domain via the Fast Fourier Transform (FFT), and subsequently discretizes the resulting frequency components through binning into variable-width intervals. Our proposed tokenizer divides the frequency range into logarithmically spaced bins that correspond to just-noticeable



(a) Frequency-based tokenizer.



(b) EnCodec-based tokenizer.

Figure 3: Two haptic tokenizers to encode vibration signals into input representations suitable for LLaMA.

differences in human frequency perception (Choi and Kuchenbecker, 2013; Israr et al., 2006). Similarly, the amplitude range is segmented into normalized levels. The tokenizer then assigns a unique token (e.g., FREQ_3_AMP_2) to each frequency-amplitude pair, encoding the signal’s spectral content into a form interpretable by LLMs.

EnCodec-based Tokenizer: EnCodec¹ is a neural audio codec that compresses audio using deep learning (Défossez et al., 2023). It consists of three main components: (1) an encoder that transforms raw audio into a lower-dimensional latent representation, (2) a quantizer that discretizes the latent features via residual vector quantization (Zeghidour et al., 2021), and (3) a decoder that reconstructs the waveform from the quantized codes. Since vibrations have perceptual and rhythmic similarities to audio signals (Bernard et al., 2022; Pätzold et al., 2023; Yun et al., 2023b; Degraen et al., 2021), we adopt the discrete codes produced by the quantizer as haptic tokens for HapticLLaMA (Figure 3(b)).

3.3.2 Input Format

After tokenization, the haptic tokens are added to the vocabulary of LLaMA tokenizer and their embeddings are updated in the learning process. Given an input tuple $\langle s, c \rangle$, where s represents the signal token sequence generated by the hap-

tic tokenizer and c denotes the category given from $\{\text{sensory, emotional, associative}\}$. For the input of HapticLLaMA, we concatenate haptic, category, textual prompt, and special tokens (e.g., $\langle \text{EOS} \rangle$) as a multimodal prompt I . We use the prompt $I = \text{“haptic signal: } \langle \text{haptic tokens} \rangle, \text{ its } \langle \text{category} \rangle \text{ description is: } \langle \text{caption} \rangle \text{.”}$. During training, we append the human-written caption as reference. The auxiliary text tokens (e.g., “its”, “description”, and “is”) can be interpreted as prompts that guide the model’s output.

3.4 Multimodal HapticLLaMA Model

Algorithm 1 outlines the two-stage training procedure of HapticLLaMA, consisting of (1) supervised fine-tuning with LoRA adaptation and (2) subsequent fine-tuning based on human feedback on generated captions, as detailed below.

3.4.1 Stage 1: Supervised Fine-Tuning

We adopt LLaMA3² as the backbone of HapticLLaMA. To align the haptic and text modalities, we incorporate haptic tokens into the LLaMA tokenizer’s vocabulary by using randomly initialized special tokens in the LLaMA vocabulary. Their embeddings are updated during training, enabling the language model to effectively interpret and utilize the haptic tokens. For efficient fine-tuning, we employ Low-Rank Adaptation (LoRA; Hu et al., 2022), which inserts trainable low-rank matrices ΔW_* into the model weights.

$$\begin{aligned} W_*^{\text{loa}} &= W_* + \Delta W_* \\ &= W_* + B_* A_*, * \in \{Q, V\} \end{aligned} \quad (1)$$

where W_*^{loa} denotes the LoRA-adapted weight matrix, and W_* refers to the original pretrained weight matrix (e.g., the query Q or value V projection weights in a transformer). The subscripts $* \in \{Q, V\}$ indicate the Query and Value projection layers, respectively, and only the matrices A_* and B_* are trainable parameters.

3.4.2 Stage 2: Fine-Tuning via RLHF

Using the trained model from Stage 1, we generate captions for each vibration signal in **VibRate** dataset (See Section 4 for details) across sensory, emotional, and associative categories and collect user ratings for the captions on 1-7 Likert scale. The rated captions are then paired into preferred

¹https://huggingface.co/facebook/encodec_24khz

²<https://huggingface.co/meta-llama/llama-3.2-3B>

Algorithm 1 HapticLLaMA Training

Require: HapticCap Train set $T = (S, D)$, VibRate signal set S' , describe category c , prompt p , LLaMA, haptic tokenizer.

Ensure: Every signal is associated with a corresponding description.

Stage 1: Supervised Fine-Tuning

- 1: Load pretrained LLaMA.
- 2: Tokenize each haptic signal into discrete haptic tokens $S = \{s_1, \dots, s_n\}$.
- 3: Include haptic tokens and update LLaMA’s vocabulary.
- 4: Formalize input based on S , p , and c into $[p_1, \dots, s_i, \dots, p_i, \dots, \langle \text{EOS} \rangle]$.
- 5: Supervised fine-tune HapticLLaMA with formalized input.
- 6: **return** initial HapticLLaMA \mathcal{H} .

Stage 2: Fine-Tuning via RLHF

- 7: Infer the caption D' of $s \in S'$ using trained \mathcal{H} .
 - 8: Construct the VibRate dataset by incorporating human ratings R .
 - 9: Construct caption preference pairs set \hat{T} based on VibRate.
 - 10: Fine-tune HapticLLaMA \mathcal{H} with \hat{T} .
 - 11: **return** The final HapticLLaMA🏆.
-

and rejected captions for fine-tuning using the DPO policy³.

Preference Pairs: To create preference pairs for DPO, we match captions rated above the midpoint of the scale with those rated below it. Rated captions are split into high-rated (positive) and low-rated (negative) groups using a 3.5 threshold. Each high-rated caption is set as the preferred choice and paired with all low-rated captions as the non-preferred (rejected) ones. The HapticLLaMA model is then trained to assign higher likelihood to the preferred response over the rejected one, aligning generation with human preferences.

4 VibRate Dataset Construction

VibRate is a diverse, manually curated dataset of 16,896 <vibration, caption, rating> tuples (see Appendix E for details on recruitment and payment). It includes 704 vibrations constructed from four diverse sources: (a) 174 vibration signals are created by varying signal parameters (e.g., frequency); (b)

180 vibrations derived from sound effect libraries by mimicking timing or applying low-pass filtering (Ternes and MacLean, 2008; Degraen et al., 2021; Yun et al., 2023b); (c) 176 vibrations generated by HapticGen (Sung et al., 2025); and (d) 174 custom-made vibrations created manually through signal transformations such as time reversal, repetition, and mixing (Schneider and MacLean, 2016; MacLean et al., 2017). For each signal, we generate four captions in each category: two using frequency-based \mathcal{H} and two from EnCodec-based \mathcal{H} . Then, we collect ratings from 44 human evaluators for the generated captions. Annotators are instructed to assign a final rating on a scale of 1 (poor) to 7 (excellent) based on two criteria: (1) the clarity and semantic accuracy of the caption, and (2) the alignment between the haptic vibration experience and caption. A higher rating indicates better quality and closer alignment with human perception. In total, 44 users participated, each evaluating captions for 32 different vibrations, yielding 16,896 <vibration, caption, human rating> samples.

5 Experiments

Datasets: We use two datasets to develop and evaluate HapticLLaMA: **HapticCap** (Hu et al., 2025) and **VibRate** on two stages respectively. **HapticCap** is the largest fully human-annotated haptic-captioned dataset, containing 92,070 haptic-text pairs, with 8-10 user-written captions per vibration, describing sensory, emotional, and associative attributes. The HapticCap dataset is divided into training, validation, and test sets (see Appendix B). We evaluate HapticLLaMA from both Stage 1 and Stage 2 on the test set.

Baselines: We set the following baselines: (1) **Random:** For each vibration in the test set, we randomly select one caption from the candidate set as the caption for that signal, (2) **Signal-agnostic:** Since LLMs can generate fluent descriptions grounded in their pretraining and guided through prompting, we examine the capability of LLaMA3.2-3B and GPT-4.5⁴ in generating captions without receiving any signal in the input (see Appendix A for details), (3) **Without LoRA Finetuning:** Fine-tuning is disabled by removing LoRA from the Frequency-based and EnCodec-based HapticLLaMA models, while keeping haptic

³https://huggingface.co/docs/trl/dpo_trainer

⁴<https://platform.openai.com/docs/models/gpt-4.5-preview>



Models	BLEU-1	BLEU-4	ROUGE-L	METEOR
Random	5.77	2.59	10.62	11.07
LLaMA (signal-agnostic)	9.39	2.71	20.67	23.84
GPT-4.5 (signal-agnostic)	11.53	4.28	26.09	28.20
Frequency tokens + LLaMA	28.46	7.48	29.49	35.24
EnCodec tokens + LLaMA	30.07	8.61	29.50	35.26
Frequency tokens + LLaMA + LoRA	40.12 \pm 5.1	23.16 \pm 4.3	42.95 \pm 3.4	50.16 \pm 2.1
EnCodec tokens + LLaMA + LoRA	41.54 \pm 3.2	24.28 \pm 2.1	44.71 \pm 2.4	54.03 \pm 2.3
Frequency HapticLLaMA  : Frequency tokens + LLaMA + LoRA + RLHF	49.43 \pm 1.4	30.86 \pm 1.5	48.74 \pm 2.1	58.95 \pm 2.1
EnCodec HapticLLaMA  : EnCodec tokens + LLaMA + LoRA + RLHF	51.36 \pm 1.3	32.06 \pm 1.2	49.56 \pm 1.6	59.98 \pm 1.8

Table 1: Results on automated n-gram metrics. Baselines include Random, LLaMA, and GPT-4.5. Ablation results show the impact of each component and the performance of HapticLLaMA with Frequency and EnCodec tokens. Green and yellow denote the best and second-best performances, respectively. \pm denotes standard deviation.

token training. (4) **Without RLHF**: Model performance with and without DPO fine-tuning from human feedback is reported.

Evaluation Metrics: Following prior work on captioning (Drossos et al., 2020; Cui et al., 2018), we evaluate haptic captions using BLEU-1 (Papineni et al., 2002), BLEU-4 (Papineni et al., 2002), ROUGE-L (Lin, 2004), and METEOR (Banerjee and Lavie, 2005) metrics. Due to the inherent ambiguity of haptic perception, each vibration in HapticCap is associated with 8–10 ground-truth reference captions. We evaluate each prediction using the reference captions and report the mean and standard deviation of the metrics to reflect performance and variability. Implementation details are provided in Appendix C.

5.1 Results

5.2 Performance on Automated Metrics

Table 1 presents the performance of our Frequency HapticLLaMA and EnCodec HapticLLaMA, the baselines, as well as the ablation results on HapticLLaMA with different haptic tokenizers, LoRA adapter, and RLHF.

The random baseline exhibits consistently poor performance across all evaluation metrics. In contrast, the signal-agnostic LLaMA and GPT-4.5 models achieve marginal improvements, attributable to their advanced language generation capabilities despite the absence of haptic signals. Ablation results demonstrate the contribution of each component in HapticLLaMA. Incorporating haptic tokens into the training vocabulary leads to consistent improvements across all evaluation metrics. LoRA fine-tuning with haptic tokenizers shows strong absolute gains between +11.47 (BLEU-1) to +18.77 (METEOR) points in all met-

rics. Applying RLHF via DPO further boosts HapticLLaMA’s performance across all metrics, including an absolute gain of +7.78 on BLEU-4 score.

EnCodec tokens consistently outperform Frequency tokens by a small margin, with EnCodec HapticLLaMA (last row) achieving the highest scores across all metrics. This result may be because the Frequency-based tokenizer emphasizes the spectral domain while neglecting temporal and rhythmic features of vibration, whereas the EnCodec-based tokenizer captures richer representations of temporal changes based on the raw vibration signal and creates efficient tokens for signal compression and reconstruction. Additionally, the length of haptic token sequences and vocabulary size produced by the EnCodec tokenizer are significantly higher than that of the Frequency-based tokenizer (see Section 5.6), suggesting that EnCodec captures more fine-grained variations in the vibration signals. These results demonstrate that both frequency and EnCodec tokens can capture input signals, and the application of DPO further enhances caption generation quality by refining the alignment between model output and human preferences.

5.3 Human Evaluation

Figure 4 presents the human evaluation results of descriptions generated by HapticLLaMA in Stage 1 and Stage 2, respectively. We use the human ratings from VibRate as the evaluation results after Stage 1. After completing both training stages, we randomly select 20 vibration signals from the HapticCap test set and generate corresponding haptic captions using the final HapticLLaMA model. Two evaluators are instructed to provide final ratings on a 1–7 scale, using the same human rating setup as in Stage 1.

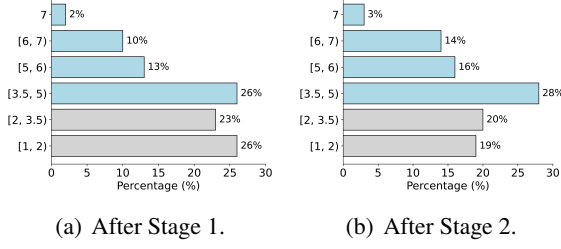


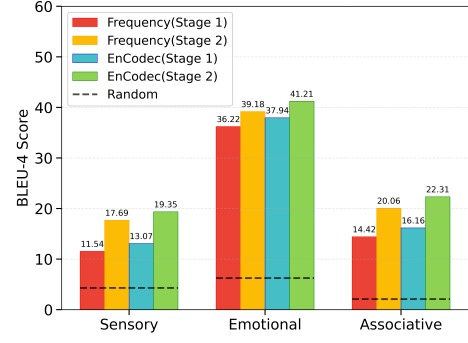
Figure 4: The distribution of human ratings after Stage 1 and 2 training for captions generated by EnCodec-based HapticLLaMA. Higher ratings are shown in blue and lower ratings are in gray.

The results show a noticeable improvement in the quality of generated captions after Stage 2, with over 61% of ratings above 3.5. Specifically, there is a visible shift in the distribution toward higher rating intervals, indicating that more captions are perceived as relevant and accurate after Stage 2. For instance, the proportion of higher ratings increases across all intervals after Stage 2, with about 10% increase in ratings above 3.5 and a particularly notable rise of about 3% and 4% in the range [5, 6) and [6, 7), respectively. In contrast, the share of lower ratings, such as [1, 2) and [2, 3.5), decreases significantly by around 7% and 3%, respectively. These results indicate that the preference alignment introduced in Stage 2 notably improves the model’s ability to generate descriptions that are semantically clearer and perceptually aligned with human judgments.

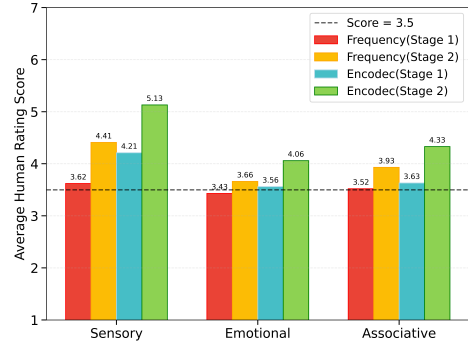
5.4 Category-Specific Performance

Figure 5 compares HapticLLaMA variants across sensory, emotional, and associative categories using (a) BLEU-4 scores (with additional METEOR results in Appendix D) and (b) average human ratings. Both illustrate that the application of DPO consistently improves performance in all categories, yielding gains in BLEU-4 and average human rating, respectively, in line with the observed improvements in overall performance.

As measured by BLEU-4, emotional captions yield the highest performance among the three categories across both HapticLLaMA variants, with the EnCodec-based model achieving the leading score of 41.2 in BLEU-4. This outcome may be attributed to the fact that emotional captions tend to show higher agreement among multiple human annotators (Hu et al., 2025), thereby facilitating more accurate and coherent generation.



(a) BLEU-4 scores.



(b) Average human ratings.

Figure 5: EnCodec-based HapticLLaMA’s performance across sensory, emotional, and associative categories.

According to human evaluations, sensory captions receive the highest ratings across all HapticLLaMA variants and stages. The discrepancy between automatic metrics and human judgments may stem from the limitations of automated methods (e.g., BLEU), which focus on surface-level textual similarity, whereas human raters are more attuned to the qualitative aspects of haptic experience. Similar discrepancy between automated metrics and human judgments are reported in image captioning literature (Elliott and Keller, 2014; Kilickaya et al., 2016).

5.5 Case Study

Figure 6 illustrates a comparison between reference ground-truth (R) and generated (G) captions for two haptic signals from the test set.

For the left signal, we observe a continuous haptic vibration featuring two prominent vibration peaks in the middle. Compared with the reference captions, the generated ones effectively capture continuity, weak intensity, and rhythm of the vibration in the sensory category. For the emotion dimension, the generated caption also successfully reflect

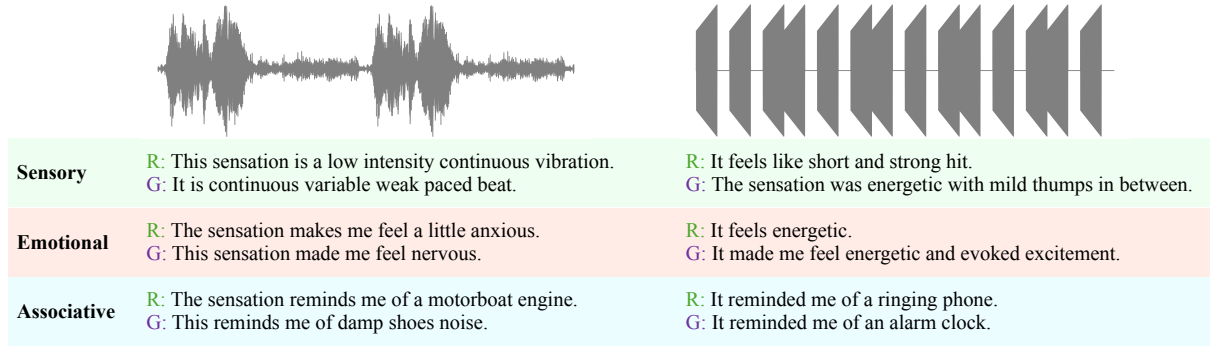


Figure 6: Case study showing two example vibrations with captions. “R” (in green) denotes reference ground-truth captions from HapticCap, while “G” (in purple) denotes captions generated by EnCodec-based HapticLLaMA.

emotional responses (e.g., “nervous” \leftrightarrow “anxious”). Regarding the associative category, HapticLLaMA broadens the range of associations from “motorboat engine” to “damp shoes”, both capturing dull, heavy sensations and demonstrating a broader and more diverse associative space. The second signal displays features that are clearly distinct from those of the first. The vibration is an intermittent sequence of pulses with pauses at regular intervals. In the sensory caption, HapticLLaMA captures the pulsing characteristic (“short hit” \leftrightarrow “mild thumps in between”) and overall energy (“strong” \leftrightarrow “energetic”). For the emotional aspect, HapticLLaMA aligns with the overall sentiment and enhances it (e.g., adding “evoked excitement”). In terms of associative categories, “ringing phone” and “alarm clock” sensations are logically similar, with both involving sequences of strong pulses in mobile phones.

5.6 Haptic Tokenizer Analysis

Table 2 presents the summary statistics of haptic tokens for two tokenizers: Frequency and EnCodec.

EnCodec-based tokenizer has a significantly larger haptic token vocabulary (1,024) compared to the Frequency-based tokenizer (278). This suggests that EnCodec can represent haptic signals with greater granularity and expressive power, consistent with the improved performance observed in the EnCodec-based HapticLLaMA. The Frequency-based tokenizer produces shorter sequences when representing a vibration signal, with an average length of 47.5 tokens, indicating a more compact representation. In contrast, EnCodec outputs a fixed-length sequence of 1,379 tokens for all signals, which may introduce redundancy but ensures uniform input sizes for downstream models. In summary, the Frequency-based tokenizer offers

Tokenizer	Item	Count
Frequency	Haptic token vocabulary size	278
	Average haptic signal length	47.5
	Max haptic signal length	52
	Min haptic signal length	12
EnCodec	Haptic token vocabulary size	1,024
	Average haptic signal length	1,379
	Max haptic signal length	1,379
	Min haptic signal length	1,379

Table 2: Summary statistics of haptic tokenizers.

higher compression efficiency, making it suitable for lightweight or real-time applications but slightly inferior performance. On the other hand, EnCodec provides more detailed and consistent representations, which may be beneficial for tasks requiring richer signal understanding.

Conclusion

We propose HapticLLaMA, the first multimodal sensory language model trained with two distinct haptic tokenizers for the haptic captioning task. HapticLLaMA formalizes haptic signals into sequences of discrete tokens through a haptic tokenizer and integrates the haptic tokens with text into a prompt template. Built on the LLaMA architecture, HapticLLaMA is trained in two stages: (1) supervised fine-tuning on the LLaMA architecture with LoRA, and (2) fine-tuning through reinforcement learning from human feedback (RLHF). Our evaluation, combining automated generation metrics (e.g., BLEU and METEOR) and human assessments, shows that HapticLLaMA can effectively perceive haptic vibration signals and provide substantial improvements over existing LLMs on the haptic captioning task. Our work advances haptic-language understanding by enabling large language models to interpret physical sensory data.

Limitations

BLEU, ROUGE, and METEOR metrics are sub-optimal for evaluating haptic captioning quality, as they primarily emphasize textual fluency and lexical overlap, while failing to account for the semantic alignment between generated captions and the underlying haptic signals. Human evaluation, although more reliable, is resource-intensive and typically conducted on a small, randomly sampled subset due to high labor costs. As a result, neither automatic nor manual evaluation methods provide a fully accurate assessment of captioning quality. HapticLLaMA focuses on vibration signals as the most accessible and diverse form of haptics, but cannot interpret other forms of haptics, such as force feedback or temperature signals.

Ethics Statement

While the idea of HapticLLaMA interpreting haptic vibration signals is appealing and holds potential for applications in human-computer interaction (HCI) and robotics, its current performance remains limited. The model achieves only around 32.06% on BLEU-4, an average human rating of 4.8 on a 7-point scale, and about 61% of captions rated above 3.5, indicating that its performance remains insufficient for real-world deployment.

Acknowledgement

We sincerely thank the volunteers for their generous contributions and invaluable efforts in providing high-quality data annotation, which has been instrumental in supporting our research. This work was supported by research grants from VILLUM FONDEN (VIL50296) and the National Science Foundation (#2339707).

References

- Jagan K Balasubramanian, Bence L Kodak, and Yasemin Vardar. 2024. Sens3: Multisensory database of finger-surface interactions and corresponding sensations. In *International Conference on Human Haptic Sensing and Touch Enabled Computer Applications*, pages 262–277. Springer.
- Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
- Sliman Bensmaïa, Mark Hollins, and Jeffrey Yau. 2005. Vibrotactile intensity and frequency information in the pacinian system: A psychophysical model. *Perception & psychophysics*, 67(5):828–841.
- Corentin Bernard, Jocelyn Monnoyer, Michaël Wiertelowski, and Sølvi Ystad. 2022. Rhythm perception is shared between audio and haptics. *Scientific Reports*, 12(1):4188.
- Ning Cheng, Changhao Guan, Jing Gao, Weihao Wang, You Li, Fandong Meng, Jie Zhou, Bin Fang, Jinan Xu, and Wenjuan Han. 2024a. Touch100k: A large-scale touch-language-vision dataset for touch-centric multimodal representation. *arXiv preprint arXiv:2406.03813*.
- Ning Cheng, You Li, Jing Gao, Bin Fang, Jinan Xu, and Wenjuan Han. 2024b. Towards comprehensive multimodal perception: Introducing the touch-language-vision dataset. *arXiv preprint arXiv:2403.09813*.
- Zebang Cheng, Zhi-Qi Cheng, Jun-Yan He, Kai Wang, Yuxiang Lin, Zheng Lian, Xiaojiang Peng, and Alexander Hauptmann. 2024c. Emotion-llama: Multimodal emotion recognition and reasoning with instruction tuning. *Advances in Neural Information Processing Systems*, 37:110805–110853.
- Inrak Choi, Yiwei Zhao, Eric J. Gonzalez, and Sean Follmer. 2021. [Augmenting Perceived Softness of Haptic Proxy Objects Through Transient Vibration and Visuo-Haptic Illusion in Virtual Reality](#). *IEEE Transactions on Visualization and Computer Graphics*, 27(12):4387–4400.
- Seungmoon Choi and Katherine J. Kuchenbecker. 2013. [Vibrotactile display: Perception, technology, and applications](#). *Proceedings of the IEEE*, 101(9):2093–2104.
- Yin Cui, Guandao Yang, Andreas Veit, Xun Huang, and Serge Belongie. 2018. Learning to evaluate image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5804–5812.
- Heather Culbertson, Juan José López Delgado, and Katherine J Kuchenbecker. 2014. One hundred data-driven haptic texture models and open-source methods for rendering on 3d objects. In *2014 IEEE haptics symposium (HAPTICS)*, pages 319–325. IEEE.
- Tor-Salve Dalsgaard, Joanna Bergström, Marianna Obrist, and Kasper Hornbæk. 2022. A user-derived mapping for mid-air haptic experiences. *International Journal of Human-Computer Studies*, 168:102920.
- Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. 2023. [High fidelity neural audio compression](#). *Trans. Mach. Learn. Res.*, 2023.
- Donald Degraen, Bruno Fruchard, Frederik Smolders, Emmanouil Potetsianakis, Seref Güngör, Antonio

- Krüger, and Jürgen Steimle. 2021. Weirding haptics: In-situ prototyping of vibrotactile feedback in virtual reality through vocalization. In *The 34th Annual ACM symposium on user interface software and technology*, pages 936–953.
- Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. 2020. Clotho: An audio captioning dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 736–740. IEEE.
- Desmond Elliott and Frank Keller. 2014. Comparing automatic evaluation measures for image description. In *52nd Annual Meeting of the Association for Computational Linguistics*, pages 452–457. Association for Computational Linguistics.
- Letian Fu, Gaurav Datta, Huang Huang, William Chung-Ho Panitch, Jaimyn Drake, Joseph Ortiz, Mustafa Mukadam, Mike Lambeta, Roberto Calandra, and Ken Goldberg. 2024. A touch, vision, and language dataset for multimodal alignment. *arXiv preprint arXiv:2402.13232*.
- Aryo Pradipta Gema, Pasquale Minervini, Luke Daines, Tom Hope, and Beatrice Alex. 2023. Parameter-efficient fine-tuning of llama for the clinical domain. *arXiv preprint arXiv:2307.03042*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.
- Guimin Hu, Daniel Hershcovich, and Hasti Seifi. 2025. [Hapticap: A multimodal dataset and task for understanding user experience of vibration haptic signals](#).
- Guimin Hu, Zirui Zhao, Lukas Heilmann, Yasemin Varadar, and Hasti Seifi. 2024. Grounding emotional descriptions to electrovibration haptic signals. *arXiv preprint arXiv:2411.02118*.
- Vladimir Iashin and Esa Rahtu. 2020. Multi-modal dense video captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 958–959.
- Ali Israr, Hong Z Tan, and Charlotte M Reed. 2006. Frequency and amplitude discrimination along the kinesthetic-cutaneous continuum in the presence of masking stimuli. *The Journal of the Acoustical society of America*, 120(5):2789–2800.
- Lynette A Jones and Nadine B Sarter. 2008. Tactile displays: Guidance for their design and application. *Human factors*, 50(1):90–111.
- Robert K Katzschnmann, Brandon Araki, and Daniela Rus. 2018. Safe local navigation for visually impaired users with a time-of-flight and haptic feedback device. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 26(3):583–593.
- Mert Kilickaya, Aykut Erdem, Nazli Ikizler-Cinbis, and Erkut Erdem. 2016. Re-evaluating automatic metrics for image captioning. *arXiv preprint arXiv:1612.07600*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Xubo Liu, Qiushi Huang, Xinhao Mei, Haohe Liu, Qiuqiang Kong, Jianyuan Sun, Shengchen Li, Tom Ko, Yu Zhang, Lilian H Tang, et al. 2022. Visually-aware audio captioning with adaptive audio-visual attention. *arXiv preprint arXiv:2210.16428*.
- Karon E MacLean, Oliver S Schneider, and Hasti Seifi. 2017. Multisensory Haptic Interactions: Understanding the Sense and Designing for It. In *The Handbook of Multimodal-Multisensor Interfaces: Foundations, User Modeling, and Common Modality Combinations*, volume 1, pages 97–142.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533.
- Ron Mokady, Amir Hertz, and Amit H Bermano. 2021. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*.
- Marianna Obrist, Sue Ann Seah, and Sriram Subramanian. 2013. Talking About Tactile Experiences. In *Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems*, pages 1659–1668.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Bastian Pätzold, Andre Rochow, Michael Schreiber, Raphael Memmesheimer, Christian Lenz, Max Schwarz, and Sven Behnke. 2023. Audio-based roughness sensing and tactile feedback for haptic perception in telepresence. In *2023 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 1387–1392. IEEE.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741.

- Oliver S Schneider and Karon E MacLean. 2016. Studying design process and example use with macaron, a web-based vibrotactile effect editor. In *IEEE Haptics Symposium (Haptics)*, pages 52–58. IEEE.
- Hasti Seifi and Karon E MacLean. 2017. Exploiting haptic facets: Users’ sensemaking schemas as a path to design and personalization of experience. *International Journal of Human-Computer Studies*, 107:38–61.
- Hasti Seifi, Kailun Zhang, and Karon E MacLean. 2015. Vibviz: Organizing, visualizing and navigating vibration libraries. In *2015 IEEE World Haptics Conference (WHC)*, pages 254–259. IEEE.
- Caitlyn E. Seim, Brandon Ritter, Thad E. Starner, Kara Flavin, Maarten G. Lansberg, and Allison M. Okamura. 2022. [Design of a Wearable Vibrotactile Stimulation Device for Individuals With Upper-Limb Hemiparesis and Spasticity](#). *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 30:1277–1287. Conference Name: IEEE Transactions on Neural Systems and Rehabilitation Engineering.
- Youjin Sung, Kevin John, Sang Ho Yoon, and Hasti Seifi. 2025. Hapticgen: Generative text-to-vibration model for streamlining haptic design. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, pages 1–24.
- David Ternes and Karon E MacLean. 2008. Designing large sets of haptic icons with rhythm. In *Haptics: Perception, Devices and Scenarios: 6th International Conference, EuroHaptics 2008 Madrid, Spain, June 10-13, 2008 Proceedings 6*, pages 199–208. Springer.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#). *CoRR*, abs/2302.13971.
- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3156–3164.
- Wenhao Wu, Haipeng Luo, Bo Fang, Jingdong Wang, and Wanli Ouyang. 2023. Cap4video: What can auxiliary captions do for text-video retrieval? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10704–10713.
- Yi Xin, Siqi Luo, Xuyang Liu, Haodi Zhou, Xinyu Cheng, Christina E Lee, Junlong Du, Haozhe Wang, MingCai Chen, Ting Liu, et al. 2024. V-petl bench: A unified visual parameter-efficient transfer learning benchmark. *Advances in neural information processing systems*, 37:80522–80535.
- Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, Xiaohui Tao, and Fu Lee Wang. 2023. Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment. *arXiv preprint arXiv:2312.12148*.
- Fengyu Yang, Chenyang Ma, Jiacheng Zhang, Jing Zhu, Wenzhen Yuan, and Andrew Owens. 2022. Touch and go: Learning from human-collected vision and touch. *arXiv preprint arXiv:2211.12498*.
- Yongjae Yoo, Taekbeom Yoo, Jihyun Kong, and Seungmoon Choi. 2015. Emotional responses of tactile icons: Effects of amplitude, frequency, duration, and envelope. In *2015 IEEE World Haptics Conference (WHC)*, pages 235–240. IEEE.
- Wenzhen Yuan, Siyuan Dong, and Edward H Adelson. 2017. Gelsight: High-resolution robot tactile sensors for estimating geometry and force. *Sensors*, 17(12):2762.
- Gyeong Yun, Minjae Mun, Jungeun Lee, Dong-Geun Kim, Hong Z Tan, and Seungmoon Choi. 2023a. Generating Real-Time, Selective, and Multimodal Haptic Effects from Sound for Gaming Experience Enhancement. In *Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems (CHI)*, CHI ’23, pages 1–17, New York, NY, USA. Association for Computing Machinery.
- Gyeong Yun, Minjae Mun, Jungeun Lee, Dong-Geun Kim, Hong Z Tan, and Seungmoon Choi. 2023b. Generating real-time, selective, and multimodal haptic effects from sound for gaming experience enhancement. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–17.
- Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. 2021. Soundstream: An end-to-end neural audio codec. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:495–507.
- Yiming Zhang, Hong Yu, Ruoyi Du, Zhanyu Ma, and Yuan Dong. 2022. Caption feature space regularization for audio captioning. *arXiv preprint arXiv:2204.08409*.

A Prompt for Signal-agnostic Experiments

In the absence of haptic signal input (signal-agnostic baseline), LLaMA is trained and tested purely on prompts “*its <category> description is: <caption>*” to evaluate its signal-agnostic capability. We provide the caption during the training and remove it from the prompt during inference. The prompts for GPT-4.5 are designed to guide the model in generating relevant haptic descriptions for sensory, emotional, and associative aspects (see Table 3), following the data collection guidelines in prior work (Hu et al., 2025). We provide several caption demonstrations in the prompt to guide the generation of GPT-4.5, as shown in Table 4.

	Role	Prompt
Sensory	system	You are experiencing a tactile sensation. Describe what you feel in one sentence.
	user	How would you describe the sensation?
Emotional	user	How does this sensation make you feel, can you attach any emotion to it?
Associative	user	Can you associate any action or object with this sensation?

Table 3: The prompt templates for GPT-4.5.

	Role	Response
Sensory	assistant	sensory: It is a steady, regular and strong pulsation
Emotional	assistant	emotional: It makes me feel calm, peaceful and focused.
Associative	assistant	associative: It’s quite like a machine gun.

Table 4: The demonstrations for GPT-4.5.

B Details of Data Split

Table 5 presents the data split of the HapticCap dataset, along with the category distributions across sensory, emotional, and associative captions.

Category	Train	Valid	Test
Sensory	24,641	2,677	3,372
Emotional	24,641	2,677	3,372
Associative	24,641	2,677	3,372
All	73,923	8,031	10,116

Table 5: The details of HapticCap.

C Implementation Details

All experiments are conducted on NVIDIA RTX A100 and RTX H100 GPUs. Due to limited computational resources, we adopt LLaMA3.2-3B as

the backbone of HapticLLaMA, and integrate Low-Rank Adaptation (LoRA) into all query and value projection layers of the Transformer architecture to enable parameter-efficient fine-tuning. In Stage 1, the batch size is set to 4, with an overall learning rate of $3e-4$. The model is trained with a generative loss. In Stage 2, we adopt the default setting of DPO trainer to optimize HapticLLaMA. The model is optimized the policy to align with human preferences through a simple classification objective, thereby improving caption quality.

D More Experimental Results

Figure 7 presents the METEOR scores of EnCodec HapticLLaMA across the three categories. The trends across the three categories are consistent with those observed in BLEU (see Figure 5(a)).

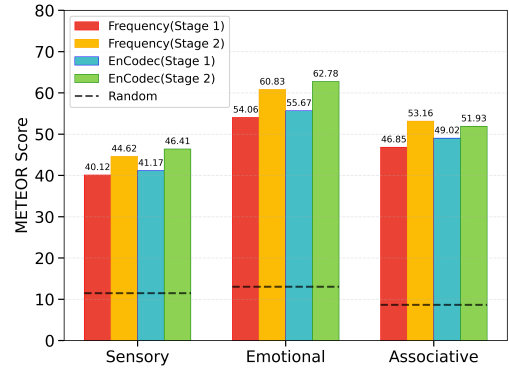


Figure 7: Performance of EnCodec-based HapticLLaMA across sensory, emotional, and associative categories, as measured by METEOR scores.

E Data Collection and Compensation

In this work, the participants were university students recruited via advertisements and snowball sampling. Each participant rated captions for 32 signals in a one-hour session and received \$15 USD in cash, which exceeds the local minimum wage.