# Context-Aware Summarization of Social Media Chat: Techniques, Challenges, and Future Directions with Large Language Models

Dr. Manuj Darbari, Priyanshu Yadav, Shwetank Rai

## I. Introduction to Chat Summarization

The proliferation of digital communication has led to an exponential increase in conversational data generated through messaging applications, social media platforms, online forums, and collaborative tools. Effectively managing and extracting value from this vast amount of information necessitates automated solutions, among which dialogue summarization plays a critical role. Dialogue summarization is formally defined as the task of condensing conversations involving two or more participants into shorter, informative versions that capture the most salient information.[1] The need for such condensation is driven by the sheer volume of dialogue data, making manual review impractical in many scenarios.[3]

Chat summarization emerges as a specialized sub-field within dialogue summarization, focusing specifically on the characteristics of informal, text-based, multi-turn conversations commonly found on social media and messaging platforms.[4] These chats often exhibit distinct linguistic features, including brevity, non-standard grammar, slang, abbreviations, and the prevalent use of emojis or emoticons.[7] The applications of dialogue and chat summarization are diverse, ranging from generating meeting minutes and summarizing customer service interactions to recapping doctor-patient consultations and digesting sprawling social media discussions.[1]

While sharing the fundamental goal of information condensation with general text summarization (e.g., summarizing news articles or scientific papers), dialogue and chat summarization present a unique set of challenges.[3] General text summarization typically deals with well-structured, single-author documents where information flow is relatively linear.[15] In contrast, dialogues are interactive, involving multiple participants whose contributions interleave and build upon each other.[1] Key information relevant to a summary might be sparsely distributed across different turns and speakers, requiring sophisticated methods to identify and synthesize it.[4] Conversations frequently exhibit topic drifts, where the focus shifts, sometimes abruptly, demanding that summarization models track these changes to capture the essence of potentially multi-threaded discussions.[1] Furthermore, dialogues are rich in interactive signals (e.g., questions, answers, confirmations, disagreements) and coreferences that are crucial for understanding the conversational flow and the importance of different utterances.[4] Spoken dialogues introduce additional complexities like disfluencies, hesitations, and speaker dynamics that must be handled.[12] The information distribution and inherent nature of dialogue, particularly informal chat, differ fundamentally from single documents. Chat summarization must contend with dynamically evolving, potentially fragmented, and often implicit information spread across participants and turns, relying heavily on shared understanding not explicitly encoded in the text.[17] Consequently, techniques optimized for document summarization often prove inadequate for the complexities of chat data, necessitating the development of specialized modeling approaches that prioritize context integration and participant interaction modeling.[4]

This report focuses on the specific challenges and opportunities within **Context-Aware Summarization of Social Media Chat**, with a particular emphasis on the application of **Bidirectional Encoder Representations from Transformers (BERT)** and other modern Natural Language Processing (NLP) techniques, including Large Language Models (LLMs) [User Query]. The emphasis on context-awareness is paramount due to the inherently fragmented, implicit, and dynamic nature of chat conversations.[17] Incorporating dialogue history, user information, or external knowledge becomes essential for accurately interpreting utterances and generating meaningful summaries. The social media environment further complicates this task by introducing unique linguistic styles, user-generated content variability, and platform-specific norms.[7] While BERT and subsequent LLMs provide powerful representational capabilities, their effective application to this specialized domain requires careful adaptation and consideration of their inherent strengths and limitations.[13] The focus on context-aware methods is therefore not merely an enhancement but a fundamental requirement for tackling the complexities

of social media chat summarization effectively.

## II. Landscape of Chat Summarization Techniques

The field of automatic text summarization, including chat summarization, has evolved significantly, broadly categorized into extractive, abstractive, and hybrid methodologies.

**Extractive, Abstractive, and Hybrid Approaches**

● **Extractive Summarization:** This approach constructs summaries by selecting and concatenating important sentences, phrases, or utterances directly from the source chat log.[4] The selection process can be based on various criteria. Early methods often relied on unsupervised statistical or heuristic features like term frequency (TF-IDF), sentence position (assuming important information appears early), or keyword matching.[16] Graph-based algorithms, such as TextRank, represent sentences as nodes and use graph centrality measures to identify key utterances.[38] Supervised methods train machine learning or deep learning models to classify utterances as summary-worthy or not, often using features derived from the text or conversation structure.[29] The primary strength of extractive methods lies in their inherent factual consistency, as they only use content present in the original source.[16] They can also be computationally simpler. However, they often suffer from poor fluency and coherence, as the selected sentences may lack smooth transitions and logical flow, especially when extracted from unstructured chats.[15] This can result in summaries that feel jarring or robotic.[15] Furthermore, extractive methods struggle to capture information that is implicitly conveyed or distributed across multiple utterances, as they rely on explicit statements within individual sentences.

● **Abstractive Summarization:** In contrast, abstractive summarization aims to generate novel sentences and phrases that capture the core meaning and salient information of the chat, potentially using words not found in the original source.[1] This approach mimics human summarization by paraphrasing and synthesizing information.[15] The goal is to produce summaries that are more concise, fluent, coherent, and human-like.[4] Historically, abstractive summarization faced significant challenges in maintaining factual accuracy (avoiding "hallucinations") and ensuring content coherence.[29] However, the advent of neural network models, particularly sequence-to-sequence architectures and large pre-trained language models, has led to substantial progress, making abstractive methods increasingly dominant.[4]

● **Hybrid Summarization:** Hybrid methods seek to combine the strengths of both extractive and abstractive approaches.[4] A common strategy involves using an extractive step to first identify key sentences or utterances, which are then fed into an abstractive model for rewriting, paraphrasing, or refinement to improve fluency and conciseness.[29] This can potentially offer a balance between the factual grounding of extraction and the readability of abstraction. Another hybrid application involves using extraction as a pre-processing step for long chats, selecting the most relevant segments to feed into an abstractive model whose input length is limited.[23]

The clear trend towards abstractive methods, accelerated by neural networks and LLMs, reflects a research emphasis on generating summaries that are not just factually grounded but also highly readable and semantically rich.[15] The advantages often cited for abstraction – fluency, coherence, human-like quality, better handling of unstructured text [15] – stand in contrast to the perceived weaknesses of extraction, such as jarring transitions and poor performance on informal dialogue.[15] The superior fluency and coherence of LLMs [29] further reinforces this direction. This suggests that many downstream applications prioritize summaries that are easily and quickly consumable, even if this introduces challenges related to maintaining strict factual accuracy.[29] Consequently, ensuring the faithfulness of abstractive summaries has become a major research focus.[1]

**Evolution: Sequence-to-Sequence and Transformer Models**

The capabilities of abstractive summarization have been significantly advanced by deep learning. Early neural approaches involved structured methods like tree or graph-based models, or generative sequence-to-sequence (Seq2Seq) models using Recurrent Neural Networks (RNNs) or Long Short-Term Memory networks (LSTMs).[29]

The introduction of the Transformer architecture [15], with its self-attention mechanism, marked a major breakthrough. Transformers allowed for better modeling of long-range dependencies and parallel computation, leading to improved performance. Building on this, pre-trained language models (PLMs) like BERT [30], BART (Bidirectional Auto-Regressive Transformer) [2], T5 (Text-to-Text Transfer Transformer) [2], PEGASUS (Pre-training with Extracted Gap-sentences for Abstractive Summarization) [38], and Longformer [13] have become the de facto standard for state-of-the-art summarization.[2] These models leverage massive unlabeled text corpora during pre-training to learn rich linguistic representations, which can then be effectively transferred to downstream tasks like summarization through fine-tuning on specific dialogue summarization datasets.[2] Models like BART and PEGASUS, with architectures specifically designed or pre-trained for generation tasks, have shown particularly strong results in abstractive summarization.[2]

**The LLM Paradigm**

More recently, the emergence of Large Language Models (LLMs) – models with billions or even trillions of parameters, such as OpenAI's GPT series, Google's PaLM/Gemini, Meta's Llama, and Mistral AI's models – represents another

significant paradigm shift.[29] LLMs exhibit remarkable capabilities due to their vast scale and training data:

- **Paradigm Flexibility:** LLMs are not constrained by predefined generative paradigms determined by specific architectures or training data, unlike earlier models.[29] They can often perform various summarization styles (extractive, abstractive) based on instructions.

- **Zero-Shot and Few-Shot Learning:** LLMs can perform tasks, including summarization, with no or very few examples provided in the prompt (in-context learning), demonstrating strong generalization abilities.[29]

- **Enhanced Quality:** LLMs often generate summaries with superior coherence, fluency, and overall quality compared to previous models, frequently aligning better with human preferences, sometimes even when scoring lower on traditional metrics like ROUGE.[29]

Leveraging LLMs for summarization typically involves one or more of the following strategies:

- **Prompt Engineering:** This involves carefully designing input prompts or instructions to guide the LLM's generation process without modifying the model's parameters.[29] The effectiveness of LLMs can be highly sensitive to the specific wording and structure of the prompt.[33] Techniques like Chain-of-Thought prompting (encouraging step-by-step reasoning) or providing detailed rubrics within the prompt might be relevant for complex summarization tasks.[23]

- **Fine-Tuning:** LLMs can be further trained (fine-tuned) on specific summarization datasets to specialize their capabilities for that particular task or domain.[2] This allows the model to adapt its internal parameters for improved accuracy and relevance.[29] Fine-tuning on data closely matching the target domain generally yields better results.[38] However, fine-tuning can be computationally expensive and may sometimes lead to a degradation of the model's general abilities (catastrophic forgetting) or require substantial labeled data.[33]

- **Knowledge Distillation:** This technique aims to transfer the summarization capabilities of a large, powerful LLM to a smaller, more efficient model.[29] This is particularly useful for deploying summarization models in resource-constrained environments where running a massive LLM is impractical.

- **In-Context Learning (ICL):** LLMs can learn to perform a task by observing a few examples (demonstrations) provided directly within the prompt, without any parameter updates.[29] This has proven effective for few-shot dialogue summarization, enabling adaptation with minimal labeled data.[44]

Interestingly, despite the dominance of abstractive methods fueled by LLMs, there is a noticeable resurgence of interest in hybrid techniques [29] and even the application of LLMs for *extractive* summarization.[32] This suggests a growing acknowledgment that pure abstraction, particularly given the known limitations of current LLMs, may not always be the ideal solution. The persistent challenges of hallucination and factual inconsistency in abstractive LLM outputs [1], coupled with difficulties in handling very long contexts effectively [14], appear to be driving this refinement. Hybrid approaches that use extraction to ground the summary or LLM-based extractive methods that leverage the model's understanding to select key sentences offer ways to harness LLM capabilities while potentially improving faithfulness and managing input length constraints.[23] This indicates a potential move towards more nuanced approaches that seek an optimal balance between the fluency offered by LLMs and the reliability associated with extraction or structured processing, especially for applications where accuracy is paramount.

**III. Unique Challenges in Social Media Chat Summarization**

Summarizing social media chat presents a distinct set of challenges that go beyond those encountered in general dialogue summarization, stemming from the unique linguistic characteristics, user behaviors, and platform dynamics inherent in these environments.

**Informality, Slang, Emojis, Abbreviations, Code-Switching**

A defining characteristic of social media chat is its pervasive informality. Conversations are often replete with slang, acronyms (e.g., LOL, BRB), non-standard spellings, grammatical errors, and creative punctuation.[8] This linguistic style deviates significantly from the more formal text found in news articles, books, or even emails, which constitute the bulk of the training data for many large language models.

Emojis and emoticons are integral components of social media communication, used to convey a wide spectrum of semantic information, including emotions, objects, actions, emphasis, tone adjustment, and topic highlighting.[7] While visually compact, their interpretation is highly context-dependent and can vary significantly based on cultural background, demographic factors, personal experience, and the specific conversational context.[7] This subjectivity poses a significant challenge for NLP models aiming for consistent interpretation; a model misinterpreting an emoji (e.g., taking sarcasm literally) could lead to inaccurate summaries.[7] Furthermore, manually annotating emojis for sentiment, intent, or meaning

across diverse contexts is costly and not scalable, hindering the development of robust emoji understanding capabilities in models.[7] While large models like ChatGPT show potential in interpreting emojis, their reliability across different contexts and potential for hallucination remain concerns.[7]

Slang presents another challenge due to its dynamic nature. New terms emerge rapidly, and usage patterns can differ significantly across various social media platforms, online communities, and age groups.[11] A summarization model trained on one set of slang terms may quickly become outdated or perform poorly when encountering new vernacular.[11] Analyzing social media slang is valuable for understanding communication trends and user behavior, but requires models capable of adapting to this linguistic evolution.[11]

Code-switching, the practice of alternating between two or more languages within a single conversation or even a single utterance, is prevalent in many multilingual communities on social media.[32] This poses a significant hurdle for summarization models, even multilingual ones, if they are not specifically designed or trained to handle mixed-language input seamlessly. Processing code-switched text requires identifying language boundaries and interpreting meaning across languages, adding substantial complexity.[32]

While some research suggests that the presence of informal language might have only a limited impact on the performance of certain NLP tasks like sentiment analysis in specific experimental setups [9], summarization arguably requires a deeper level of semantic understanding. The cumulative effect of pervasive informality, evolving slang, ambiguous emojis, and potential code-switching likely presents a significant obstacle to accurately extracting and synthesizing the core content of social media chats. Related phenomena, such as sarcasm, which is often conveyed through subtle informal cues and context, are known to be particularly challenging for NLP models to detect and interpret correctly.[9] Techniques related to formality transfer, aiming to convert informal text to a more formal style, represent a related research direction that could potentially be leveraged in pre-processing for summarization.[56]

**Topic Shifts and Multi-Party Dynamics**

Unlike structured documents, social media chats often lack clear organizational markers like paragraphs or headings.[12] Conversations can meander, with topics shifting frequently and unpredictably, sometimes within a short span of turns.[1] Identifying coherent conversational segments related to specific topics and tracking the flow of salient information across these shifts is a significant challenge.[12] Models need to discern when a topic changes and determine which information from different segments is relevant to the overall summary. Feeding long, multi-topic dialogues directly into a summarization model can dilute its focus and hinder its ability to identify key utterances, potentially leading to incomplete or fragmented summaries.[3]

The multi-party nature of many social media chats (e.g., group chats, comment threads) adds further complexity.[1] Information pertinent to a summary might be contributed piecemeal by different participants across various turns. The summarization model must effectively aggregate these scattered pieces of information, resolve coreferences (e.g., understanding who "he" or "they" refers to across different speakers' turns), and synthesize a coherent narrative.[4] Understanding the roles, intentions, and perspectives of different participants can be crucial for generating a comprehensive summary, but modeling these dynamics is inherently difficult.[1]

**Noise, User-Generated Content Issues, and Data Sparsity**

Social media platforms are inherently noisy environments.[24] Chat logs can be cluttered with irrelevant messages, off-topic discussions, interruptions, incomplete sentences, typographical errors, repetitions, and automated messages (e.g., bot responses).[24] Summarization models must be robust to this noise, capable of filtering out extraneous information and focusing on the substantive content. The sheer volume and "always-on" nature of social media content make manual curation infeasible and necessitate automated, noise-tolerant processing techniques.[24]

The quality of User-Generated Content (UGC) varies dramatically.[26] While users often value authenticity and raw, unfiltered content [25], this authenticity often manifests as the noisy and informal characteristics discussed earlier. Summarization algorithms face the task of discerning valuable insights within this heterogeneous content stream, effectively ignoring low-quality or irrelevant contributions while capturing genuine information.[24]

Data availability presents another hurdle. While massive datasets of raw chat conversations exist (e.g., WildChat [57], LMSYS-Chat-1M [58]), obtaining high-quality reference summaries for these conversations, particularly for the diverse range of social media platforms, is challenging.[2] Creating reliable summaries requires significant human effort and cost, and privacy concerns associated with personal conversations further complicate data collection and sharing.[2] This relative scarcity of labeled summary data necessitates the use of techniques like data augmentation, unsupervised methods, or few-shot learning approaches that can function effectively with limited supervised examples.[4]

Furthermore, summarizing UGC raises significant ethical considerations. Fairness is a key concern, as summarization

algorithms, in their process of selecting salient information, might inadvertently under-represent the viewpoints of certain demographic or social groups present in the conversation, especially if these groups use non-standard language patterns that the model struggles with.[27] Toxicity is another major issue; social media content can contain harmful, offensive, or discriminatory language.[59] Summarization models must be designed to handle such content appropriately, avoiding the amplification or legitimization of toxic messages within the summary.[59]

The combination of these challenges – informality, topic shifts, multi-party dynamics, noise, and ethical considerations – creates a uniquely complex environment for summarization. The difficulties are not merely additive; they interact synergistically. For instance, interpreting an informal slang term might require understanding the preceding turns from multiple speakers (multi-party dynamics) within the current micro-topic (topic shift), all while filtering out surrounding noise. An emoji's intended meaning could hinge on context established many turns earlier, demanding robust long-range dependency modeling.[7] Given that models often struggle with utilizing long context effectively [47], this interplay makes robust social media chat summarization exceptionally difficult. It necessitates models capable of handling informality *while simultaneously* tracking multiple speakers *across* topic shifts *within* noisy data streams *over* potentially extensive conversational histories. Addressing these challenges likely requires holistic approaches that consider these interactions, rather than attempting to solve each problem in isolation. Context-awareness, therefore, becomes even more indispensable in navigating this synergistic complexity.

Moreover, the subjective nature of emoji interpretation [7] and the rapid evolution of slang [11] suggest that static models trained on fixed datasets are likely to degrade in performance over time or fail when applied to different user communities or platforms. A model proficient in understanding the language used on Twitter yesterday might struggle with the nuances of a specific Facebook group today. This points towards a need for models capable of continuous learning or rapid adaptation, perhaps through ongoing fine-tuning, few-shot learning tailored to new communities, or architectures inherently designed to handle linguistic dynamism and subjective interpretations.

Finally, the presence of noise, the imperative for fairness [27], and the need to manage toxic content [59] introduce critical ethical dimensions alongside the technical hurdles. The process of summarization inherently involves selection and emphasis. If not designed with care, models could inadvertently amplify harmful narratives present in the UGC or erase the perspectives of minority groups, leading to biased or harmful outputs. Therefore, developing effective social media chat summarizers requires not only technical sophistication but also explicit consideration of fairness, bias mitigation, privacy preservation, and potentially the integration of toxicity detection and filtering mechanisms. This adds significant layers of complexity to the research and development process.

## IV. Context-Awareness in Dialogue Summarization

Given the nature of dialogue, especially informal chat, context is not merely helpful but essential for accurate understanding and effective summarization. Context-awareness refers to the ability of a summarization system to leverage information beyond the immediate utterance being processed. This context can encompass various sources, including the preceding conversation history, information about the participants, and relevant external knowledge.

### Leveraging Dialogue History

The most fundamental form of context is the dialogue history itself. Past utterances provide the necessary background to interpret context-dependent phrases (e.g., "What about that one?"), resolve ambiguities like pronoun references ("She said it was fine"), understand implicit meanings, and maintain coherence across turns.[17] A system summarizing a chat turn-by-turn without access to prior conversation would likely produce fragmented and inaccurate results.

A primary challenge in leveraging dialogue history is its potential length. Many conversations, particularly in meetings or extended chat sessions, can span hundreds or thousands of turns, exceeding the typical input limits of standard Transformer-based models due to their quadratic computational complexity with respect to input sequence length.[14] Various strategies have been developed to address this:

● **Full History Input:** For shorter conversations that fit within the model's context window, the entire history can be provided as input.

● **History Selection/Filtering:** Identifying and selecting only the most relevant previous turns to include in the context, potentially based on recency, topic similarity, or explicit relevance modeling.[19] Active research explores methods to explicitly denoise the history or assess the relevance and usefulness of historical turns for the current summarization task.[19]

● **History Summarization/Compression:** Generating a concise summary of the earlier parts of the conversation and using this summary as context when processing later parts.[21] This allows information from the

distant past to be retained in a compressed form.

● **Hierarchical or Segmented Processing:** Breaking the long dialogue into smaller segments or chunks, processing each segment (potentially incorporating context from previous segments), and then aggregating the results.[4]

● **Retrieval-Augmented Generation (RAG):** Using retrieval mechanisms to find the most relevant snippets from the long history based on the current turn or a specific query, and then feeding only these retrieved snippets along with the current input to the summarization model.[20]

## Incorporating User Information and Personalization

Traditional dialogue summarization approaches often generate generic summaries, overlooking the specific interests, goals, or background of the user who might consume the summary, or the differing perspectives of the participants involved.[2] Context-aware systems offer the potential for personalization by incorporating user-specific information.

This user information could include explicit user profiles, preferences inferred from past interactions or dialogue history, or the user's role in the conversation.[2] For example, a summary of a product discussion chat might highlight different aspects depending on whether the user is a potential buyer, a support agent, or a product developer. Tailored summaries are likely to be perceived as more relevant and useful, potentially enhancing user engagement with conversational systems.[20]

Query-based summarization represents one mechanism for incorporating user interest, where the summary is generated specifically to answer a user's question about the dialogue.[2] Developing personalized summarization systems requires either datasets specifically annotated for this purpose or methods to automatically generate relevant training data, such as synthesizing query-dialogue-summary triples, potentially leveraging LLMs themselves for this data creation process.[2]

## Utilizing External Knowledge

Conversations frequently reference real-world entities, events, or concepts that may require external knowledge beyond the dialogue itself for full comprehension. For instance, summarizing a chat discussing a specific technical problem might benefit from accessing a relevant technical manual, or summarizing a discussion about a current event might require background information about that event.

Integrating external knowledge sources – such as knowledge graphs (KGs), databases, domain-specific documents, or even the broader web – can significantly enhance a summarization model's understanding of the dialogue content and enable the generation of more informative and accurate summaries.[22] Techniques for knowledge integration often involve a retrieval step, where relevant external facts or documents are identified based on the dialogue context, followed by an encoding step where this retrieved knowledge is fused with the dialogue representation before being fed to the summarizer.[22] Specialized architectures like Graph Attention Networks (GATs) can be employed to effectively assimilate structured knowledge (like KGs) with the textual dialogue data, capturing dependencies between the conversation and external facts.[22] The RAG paradigm is also highly relevant here, retrieving external knowledge instead of, or in addition to, dialogue history snippets.[54]

## Query-based and Goal-Oriented Summarization

Moving beyond generic summaries, context-awareness enables more targeted forms of summarization driven by specific user needs or predefined goals.

● **Query-based Summarization:** As mentioned, this approach generates a summary specifically tailored to answer a user's query about the conversation.[2] The query provides explicit context regarding the user's information need. Datasets like QMSum are designed to support research in this area, providing query-dialogue-summary triples for meeting scenarios.[14]

● **Goal-Oriented/Rubric-Driven Summarization:** This aims to produce summaries that fulfill a specific function or adhere to predefined criteria (rubrics) based on the application context.[23] For example, a summary of a customer support chat might be required to capture the specific problem reported, the troubleshooting steps taken, and the final resolution, while a summary of a brainstorming session might focus on capturing novel ideas generated. These rubrics or goals can be specified, potentially in natural language, to guide the summarization process, for instance, during the fine-tuning of an LLM.[23] Related research in intent-aware dialogue systems, which model the underlying intentions or goals driving the conversation (e.g., using Hidden Markov Models combined with LLMs), could provide valuable techniques for informing goal-oriented summarization.[62]

The pursuit of context-awareness reveals its multi-faceted nature. It extends beyond merely considering past utterances to potentially encompass user identity, preferences, roles, and external world knowledge.[2] Integrating these diverse and heterogeneous sources (sequential text, structured profiles, knowledge graphs, unstructured documents) effectively within

a single summarization framework presents a significant architectural challenge. Simple concatenation is unlikely to suffice, especially given the scale of potential context. This points towards the need for hybrid architectures that combine retrieval mechanisms [20], specialized encoders (e.g., graph networks for KGs [22]), and sophisticated sequence models. The mechanism for fusing these disparate information streams remains a key area for research.

Furthermore, the increasing interest in query-based [2] and goal-oriented or rubric-driven summarization [23] signals a potential evolution in the definition of the task itself. It suggests a shift away from generating *a single, universally optimal* summary towards generating *the most appropriate* summary given a specific user need, query, or predefined goal. This positions dialogue summarization closer to controllable text generation and question answering, where the output is explicitly conditioned on external constraints or instructions.[2] This evolution has profound implications for model design (requiring controllability), training data (requiring context/query/goal annotations), and evaluation (requiring metrics that assess goal fulfillment).

However, the effective utilization of dialogue history, a cornerstone of context-awareness, is directly challenged by known limitations of current LLMs. The "lost-in-the-middle" phenomenon, where models exhibit weaker awareness of information situated in the middle of long input contexts compared to the beginning or end, is particularly concerning for dialogue summarization.[47] Social media chats, often being lengthy and meandering, may contain crucial details or context shifts buried deep within the conversation history. Naively providing the entire history as input to an LLM might result in summaries that overlook this critical mid-conversation information. Therefore, research in context-aware summarization must actively confront this LLM limitation. Strategies beyond simple concatenation are essential, including explicit relevance assessment of historical turns [19], hierarchical processing [4], retrieval pipelines [14], or the development of novel methods and architectures specifically designed to enhance context utilization across the entire input sequence.[55]

## V. BERT and Advanced NLP Techniques for Chat Summarization

The advent of Transformer-based pre-trained language models (PLMs), starting with BERT and evolving into today's sophisticated LLMs, has profoundly impacted chat summarization research and practice.

**Application and Performance Analysis (Strengths/Limitations)**

BERT and its contemporaries (e.g., RoBERTa, ELECTRA) and subsequent sequence-to-sequence PLMs like BART, T5, and PEGASUS have become foundational components in many state-of-the-art summarization systems.[4] The dominant paradigm involves pre-training these models on massive text corpora and then fine-tuning them on specific dialogue summarization datasets.[4] Specialized adaptations like BERTSUM were developed to apply BERT effectively to extractive summarization by adding summarization-specific layers.[30]

The emergence of LLMs (e.g., GPT-3/4, Llama, Mistral, PaLM, Gemini) has further advanced the field, demonstrating impressive zero-shot and few-shot summarization capabilities.[29] These models often generate summaries perceived by humans as highly fluent, coherent, and relevant, sometimes surpassing fine-tuned models in human evaluations even if not in automated metrics.[47] Task-specific LLMs, such as those fine-tuned for particular domains like mental health counseling (e.g., MentalLlama), also show considerable promise.[13]

The **strengths** of these models for chat summarization include:

● **Deep Language Understanding:** Pre-training enables them to capture complex linguistic patterns, semantic relationships, and contextual nuances.[31]

● **Semantic Capture:** They can go beyond surface-level features to understand the underlying meaning of the conversation.

● **Fluency and Coherence:** Especially LLMs, they excel at generating natural-sounding and logically flowing text.[29]

● **Transfer Learning:** Pre-trained knowledge significantly reduces the amount of task-specific data needed for effective performance compared to training models from scratch.[31]

● **Flexibility (LLMs):** LLMs can often adapt to different summarization requirements (e.g., length, focus) through prompting.[29]

However, these models also exhibit significant **limitations** in the context of chat summarization:

● **Context Length Constraints:** The self-attention mechanism in standard Transformers has computational complexity quadratic in the input sequence length ($O(n2)$), making it prohibitively expensive to process very long chat histories directly.[14] While models with more efficient attention mechanisms (e.g., Longformer's sparse attention [13]) or LLMs supporting extended context windows (e.g., up to 128k tokens via techniques like positional interpolation [47]) have been developed [13], a critical issue remains.

- **Ineffective Context Utilization:** Evidence suggests that even models capable of processing long inputs do not utilize the entire context window effectively. They often exhibit a "lost-in-the-middle" phenomenon, paying more attention to information at the beginning and end of the context while potentially ignoring crucial details in the middle.[47] This is particularly problematic for summarizing long, evolving chats.

- **Faithfulness and Hallucination:** Abstractive models, and LLMs in particular, are prone to generating content that is not factually supported by the source text (hallucination).[1] This can range from minor inaccuracies to entirely fabricated statements. In dialogue summarization, a specific challenge is the generation of "circumstantial inferences" – statements that seem plausible given the conversational context but lack direct textual evidence.[51] Ensuring faithfulness is a major ongoing research challenge.

- **Prompt Sensitivity (LLMs):** The performance of LLMs in zero-shot or few-shot settings can be highly dependent on the exact phrasing and structure of the input prompt or instruction.[29] Finding optimal prompts can require significant experimentation, and performance can drop steeply with suboptimal prompts.[33]

- **Computational Cost:** Training and deploying large PLMs and especially LLMs require substantial computational resources (GPUs, memory), limiting accessibility and increasing operational costs.[29]

- **Domain Adaptation Challenges:** Models pre-trained primarily on formal text (like web crawls or books) may struggle to handle the unique linguistic features of informal social media chat (slang, emojis, noise) without specific adaptation through fine-tuning or tailored prompting.[33]

- **Evaluation Difficulties:** As discussed previously, standard automatic metrics like ROUGE often fail to capture the nuances of summary quality, especially for fluent LLM outputs, and evaluating faithfulness remains a major hurdle.[45]

The tension between the remarkable fluency and coherence of LLM-generated summaries (which often leads to higher human preference ratings [47]) and their susceptibility to factual errors or hallucinations [48] is a central theme. This tension appears particularly acute in dialogue summarization. While studies in news summarization have sometimes found LLMs produce fewer inconsistencies than older models [51], recent work focusing on dialogue summarization revealed significant inconsistency rates (over 30% in one study) even in summaries generated by powerful models like GPT-4.[51] This discrepancy might stem from the inherent nature of conversation, which is often less explicitly factual and more reliant on inference and shared understanding compared to news reports. LLMs, optimized for coherent generation, might readily generate plausible "circumstantial inferences" to fill perceived gaps in the conversational narrative, making faithfulness particularly challenging to achieve and evaluate in the chat domain.[51]

**Fine-tuning Strategies and Prompt Engineering for Chat Data**

Adapting powerful PLMs and LLMs to the specific characteristics of chat summarization typically involves fine-tuning or prompt engineering.

- **Fine-tuning:** This remains a prevalent approach, where a pre-trained model (e.g., BART, T5, Llama) is further trained on a dataset of chat-summary pairs, such as SAMSum or DialogSum.[2] The key advantage is specialization; the model learns the specific patterns and style of the target data. Fine-tuning on datasets from a similar domain (e.g., chit-chat for social media summarization) is crucial for optimal performance.[38] Instruction tuning, a form of fine-tuning where models learn to follow natural language instructions for various tasks, has been shown to improve summarization performance and generalization compared to models only pre-trained on raw text.[33] Hybrid strategies might involve fine-tuning an abstractive model using summaries initially generated by an extractive phase, which can be beneficial when high-quality human reference summaries are scarce.[30]

- **Prompt Engineering:** For leveraging the zero-shot or few-shot capabilities of LLMs, careful prompt design is critical.[29] This includes formulating clear and effective instructions (e.g., "Summarize this chat conversation focusing on the key decisions made") [33] and potentially providing illustrative examples within the prompt (In-Context Learning, ICL).[29] For goal-oriented summarization, prompts can incorporate specific rubrics or criteria the summary should meet.[23] Given the sensitivity of LLMs to prompts [33], techniques for automatically optimizing prompts (Automatic Prompt Engineering, APE [54]) could be valuable, although their application to the nuances of chat summarization needs exploration.

The high sensitivity of LLMs to input prompts [33] poses a practical challenge for summarizing diverse and noisy social media chat. A prompt meticulously crafted for summarizing a formal customer service interaction might yield poor results when applied to an informal, slang-filled chat among friends. Manually designing and selecting optimal prompts for every conceivable chat scenario is infeasible. This suggests that relying solely on zero-shot prompting might lead to fragile and inconsistent performance in real-world social media applications. This points towards a need for more robust approaches,

potentially involving methods that learn to adapt prompts based on the input chat's characteristics, meta-learning for prompt generation, or prioritizing fine-tuning strategies that embed robustness within the model's parameters, making it less dependent on specific prompt phrasing.

### Addressing Long Context Challenges in Conversations

Effectively handling the long sequences inherent in many chat conversations is critical for context-aware summarization. Several strategies are employed:

● **Specialized Architectures:** Models like Longformer utilize efficient attention mechanisms (e.g., combining local windowed attention with task-motivated global attention) to scale to longer inputs with near-linear complexity, rather than quadratic.[13] Architectural choices regarding positional encoding (e.g., absolute vs. relative embeddings, Rotary Position Embeddings (RoPE)) also influence how models handle long sequences.[54]

● **Input Segmentation / Chunking:** A common practical approach is to divide the long dialogue into smaller, overlapping or non-overlapping chunks that fit within the model's context window.[23] Summaries can be generated for each chunk independently and then concatenated, or generated iteratively, where the summary of one chunk is passed as context along with the next chunk.[23]

● **Hierarchical Modeling:** These models process the dialogue at multiple levels of granularity. For example, encoding individual utterances, then aggregating utterance representations into segment representations, and finally generating a summary based on the segment-level information.[4]

● **Retrieval-based Methods (Retrieve-then-Summarize):** This popular two-stage pipeline first employs a retrieval component to identify the most relevant utterances or passages from the long dialogue (based on heuristics, learned models, or a specific query), and then feeds only these retrieved parts to a standard summarization model.[4] This approach has demonstrated strong performance on long dialogue summarization benchmarks [14] and aligns with the RAG framework.[54]

● **Context Compression/Summarization:** Instead of discarding past context, earlier parts of the conversation can be summarized recursively, with the condensed summary serving as context for processing subsequent parts.[21]

● **Improving Intrinsic Context Awareness:** Research is ongoing to directly improve how well LLMs utilize information across their entire context window. This includes investigating the causes of the "lost-in-the-middle" effect [47] and developing techniques to mitigate it, such as novel attention patterns or methods like Mixture of In-Context Experts (MoICE) proposed for models using RoPE.[55]

The challenge of long context processing underscores a crucial point: simply increasing the maximum token limit of LLMs (e.g., to 100k or more [47]) does not automatically guarantee effective summarization of long chats. The evidence suggesting poor utilization of information in the middle segments [47] implies that architectural innovations focusing on *how* context is attended to and utilized (e.g., efficient attention, improved positional encoding) or strategic input structuring (e.g., retrieval pipelines, hierarchical processing) are likely more critical than merely scaling the raw window size. For chat summarization, where essential context or topic conclusions might be buried deep within a lengthy history, ensuring that the model can access and leverage this information effectively is paramount. Currently, retrieval-based or hierarchical methods might offer more practical and reliable solutions than relying solely on end-to-end processing with extremely long-context LLMs whose internal utilization patterns are not yet fully understood or optimized.

### VI. Datasets and Evaluation

Progress in chat summarization research heavily relies on the availability of suitable datasets for training and evaluation, as well as appropriate metrics for assessing model performance.

### Relevant Public Datasets

Developing and evaluating chat summarization models requires datasets specifically designed for dialogue, as standard text summarization corpora like CNN/Daily Mail or XSum do not adequately capture the interactive, multi-participant, and often informal nature of conversations.[4] Several datasets have been created for dialogue and chat summarization research:

● **Chat/Dialogue Focused:**
○ **SAMSum Corpus:** Contains over 16,000 messenger-like chat dialogues (chit-chat style) with high-quality, human-written abstractive summaries.[6] Its focus on informal, multi-turn conversations makes it highly relevant for social media chat summarization research.
○ **DialogSum:** Comprises 13,460 dialogues sourced from various real-life scenarios, including daily

conversations (Dailydialog), task-oriented dialogues (DREAM, MuTual), and an English speaking practice website, paired with human-written abstractive summaries and topic labels.[33] Its diversity makes it another valuable resource.

○ **WildChat / LMSYS-Chat-1M:** These are large-scale datasets containing millions of real-world user interactions with LLMs like ChatGPT.[57] They offer unparalleled scale and diversity in terms of prompts, topics, and languages. However, they typically lack curated reference summaries, making them more suitable for pre-training, studying user behavior, or unsupervised summarization approaches, rather than direct supervised evaluation of summarization quality unless summaries are generated or annotated separately. Access often requires agreeing to specific terms.[58]

○ **Webis-TLDR-17 Corpus:** A large dataset (~3 million pairs) mined from Reddit, consisting of posts and user-provided "Too Long; Didn't Read" (TLDR) summaries.[6] This is directly relevant to social media summarization, capturing user-generated summaries in a specific online community context.

● **Meeting Focused (Relevant for Long Conversation Aspects):**

○ **AMI Corpus / ICSI Meeting Corpus:** Widely used benchmarks containing transcripts and annotations (including abstractive summaries for AMI) of multi-party meetings.[4] Useful for studying summarization of long, structured conversations.

○ **QMSum:** A query-based meeting summarization dataset across multiple domains, where summaries are generated in response to specific questions about the meeting content.[4] Valuable for research on controllable and long-dialogue summarization.

● **Other Domains:** Datasets exist for customer service (e.g., CSDS, TODSum [4]), email threads (e.g., EMAILSUM [4]), medical dialogues [1], TV/movie scripts (e.g., SummScreen [4]), and more. While potentially useful for studying specific dialogue phenomena (e.g., role modeling in customer service, formality in email), their direct relevance to informal social media chat may be limited.

Despite the availability of datasets like SAMSum, DialogSum, and Webis-TLDR-17, a gap might still exist in capturing the full spectrum of social media interactions. Platforms like Twitter (with its unique threading, brevity, and public nature), Instagram (integrating visual and textual content), Facebook (group discussions, comment chains), and ephemeral messaging apps present distinct characteristics that may not be fully represented in current benchmark datasets. Research findings validated solely on existing datasets might, therefore, face challenges when generalizing to the diverse realities of live social media platforms. Webis-TLDR-17 is a notable step towards platform-specific data (Reddit).[6] There remains a need for either more diverse datasets covering various platforms or research explicitly focused on cross-platform generalization and adaptation.

**Table 1: Overview of Key Dialogue/Chat Summarization Datasets**

| Dataset Name | Domain | Size (Dialogues/Pairs) | Annotation | Availability | Social Media Relevance |
|---|---|---|---|---|---|
| SAMSum [6] | Chat (Messenger-like) | ~16,000 pairs | Human, Abstractive | Public (e.g., Hugging Face Datasets) | High |
| DialogSum [68] | Dialogue (Mixed) | ~14,000 pairs | Human, Abstractive, Topic Labels | Public (e.g., Kaggle, Hugging Face) | High |
| Webis-TLDR-17 [6] | Social Media (Reddit) | ~3,000,000 pairs | User-generated (TLDRs), Abstractive | Public (Webis) | High (Reddit specific) |
| QMSum [61] | Meeting (Multi-domain) | ~1,800 transcripts | Human, Query-based Abstractive, Extractive | Public (GitHub) | Low (Long context) |
| AMI Corpus [4] | Meeting | ~100 hours | Human, Abstractive, | Public (AMI Corpus website) | Low (Long context) |

| | | | Extractive, etc. | | |
|---|---|---|---|---|---|
| WildChat [57] | LLM Chat (Real-world) | ~1,000,000 convos | None (Raw dialogues) | Public (Allen AI, requires consent) | Medium (User prompts) |
| LMSYS-Chat-1M [58] | LLM Chat (Real-world) | ~1,000,000 convos | None (Raw dialogues) | Public (LMSYS, requires agreement/request) | Medium (User prompts) |

*(Note: Size figures are approximate. Availability may be subject to licenses or specific access procedures.)*

This table provides a comparative overview, aiding researchers in selecting appropriate datasets based on domain, annotation type, scale, and relevance to social media chat. It highlights the trade-off between datasets with high-quality summaries (SAMSum, DialogSum) and larger datasets of raw conversations (WildChat, LMSYS-Chat-1M) or user-generated summaries (Webis-TLDR-17).

**Standard Evaluation Metrics**

Evaluating the quality of generated summaries is crucial for measuring progress. Several automatic metrics are commonly used:

- **ROUGE (Recall-Oriented Understudy for Gisting Evaluation):** This family of metrics remains the most widely adopted standard.[1] It works by comparing the generated summary against one or more human-written reference summaries based on overlapping units:
  - *ROUGE-N* (e.g., ROUGE-1, ROUGE-2) measures the overlap of n-grams (unigrams, bigrams). ROUGE-1 relates to content coverage, while ROUGE-2 is sometimes associated with fluency.[4]
  - *ROUGE-L* measures the longest common subsequence (LCS) between the generated and reference summaries, reflecting structural similarity and sentence-level ordering.[4] ROUGE primarily focuses on recall – how much of the reference summary is captured in the generated one.[72]

- **BERTScore:** Introduced to address the limitations of lexical overlap metrics like ROUGE, BERTScore leverages contextual embeddings from models like BERT to compute the similarity between tokens in the generated and reference summaries.[13] By comparing embeddings using cosine similarity, it aims to capture semantic similarity even when different wording is used.[63] It can report precision, recall, and F1 scores [66] and generally shows better correlation with human judgments than ROUGE, especially for generative tasks.[63]

- **BLEU (Bilingual Evaluation Understudy):** Originally developed for machine translation, BLEU measures n-gram precision (how many n-grams in the generated text appear in the reference) and includes a brevity penalty to discourage overly short outputs.[38] While sometimes used in summarization, its focus on precision makes it conceptually different from recall-oriented ROUGE. Some studies suggest its scores might correlate with ROUGE scores in certain contexts.[38]

- **METEOR (Metric for Evaluation of Translation with Explicit ORdering):** Another metric primarily from translation, METEOR performs unigram matching based on exact words, stems, and synonyms, incorporating a fragmentation penalty based on word ordering.[71] It aims to balance precision and recall.[72]

**Challenges in Evaluation**

Despite the availability of these metrics, evaluating chat summarization, particularly abstractive summaries generated by LLMs, remains a significant challenge:

- **Reference Dependence and Cost:** Metrics like ROUGE and BERTScore fundamentally rely on comparing the generated summary to one or more human-written reference summaries.[63] Creating these gold-standard references is labor-intensive and expensive, limiting the scale and diversity of evaluation datasets.[2] Furthermore, a single reference may not capture all possible valid ways to summarize a conversation.

- **Poor Correlation with Human Judgment:** A major limitation, especially for ROUGE, is its often low correlation with human assessments of overall summary quality, including aspects like fluency, coherence, and factual consistency.[45] Models can achieve high ROUGE scores by simply extracting keywords while producing incoherent text, or conversely, fluent and accurate abstractive summaries might be penalized for using different wording than the reference.[63] LLM-generated summaries, often preferred by humans for their fluency, may paradoxically score lower on ROUGE than summaries from older models.[47]

- **Inability to Assess Faithfulness:** Standard overlap-based metrics (ROUGE, BLEU, METEOR) and even semantic similarity metrics like BERTScore do not reliably measure factual consistency or detect hallucinations.[48] A summary might have high overlap with a reference but still contain factual errors, or conversely, be factually

accurate but use novel phrasing, leading to a low score. Evaluating faithfulness is critical, especially for abstractive LLMs, but typically requires meticulous human verification or specialized automated techniques.[48] Even human evaluation for faithfulness can be challenging and inconsistent, with annotators sometimes overlooking subtle errors in fluent summaries.[50]

● **Rise of LLM-based Evaluators:** An emerging approach uses LLMs themselves as evaluators (often termed "LLMs-as-Judges").[47] By providing the LLM with the source chat, the generated summary, and specific instructions (e.g., "Rate the faithfulness of this summary on a scale of 1-5" or "Identify any factual inconsistencies"), these models can provide automated assessments of qualities beyond lexical overlap. However, this approach is still under development and faces its own challenges. LLMs can be biased, sensitive to prompt wording, and sometimes "fooled" by the very fluency they excel at generating, potentially overlooking factual errors in well-written but inaccurate summaries.[48] Research is exploring ways to improve their reliability, for instance, through multi-agent debate frameworks where different LLM instances argue for opposing stances (e.g., faithful vs. unfaithful) to reach a more robust judgment.[48]

● **Need for Task-Specific Metrics:** Evaluating context-aware, personalized, or goal-oriented summaries may require metrics that go beyond generic quality assessment. Metrics might need to explicitly measure how well the summary incorporates the required context, reflects user preferences, or fulfills the specified goal or rubric.

**Table 2: Comparison of Evaluation Metrics for Chat Summarization**

| Metric | Focus | Pros | Cons/Limitations (esp. for Chat Summarization) |
|---|---|---|---|
| ROUGE-N/L [71] | Lexical Overlap (n-grams, LCS), Recall | Simple, Fast, Widely used standard | Low correlation with human judgment (esp. fluency, coherence, faithfulness), sensitive to wording, requires references, struggles with informality |
| BERTScore [71] | Semantic Similarity (Contextual Embeddings) | Better correlation with human judgment than ROUGE, captures semantics | Still requires references, may not fully capture faithfulness or long-range coherence, computational cost higher than ROUGE |
| BLEU [72] | Lexical Overlap (n-grams), Precision, Brevity Penalty | Fast, common in generation tasks | Primarily for translation, precision focus less suitable for summarization recall goal, similar lexical limitations to ROUGE |
| METEOR [71] | Word Overlap (incl. synonyms, stems), Precision/Recall Balance | Considers synonyms/stems, balances P/R | Primarily for translation, requires linguistic resources, less common in summarization literature |
| Human Eval [50] | Fluency, Coherence, Faithfulness, Relevance, Overall Quality | Gold standard for quality assessment, captures nuances | Slow, Expensive, Subjective, Potentially inconsistent, difficult to scale |
| LLM-as-Judge [74] | Customizable (Fluency, Faithfulness, Coherence, etc.) | Scalable, potentially captures nuances beyond overlap, no reference needed (sometimes) | Reliability concerns (fooled by fluency), potential bias, sensitivity to prompts, consistency |

| | | | issues, evaluation methodology evolving |
|---|---|---|---|

This table highlights the trade-offs involved in choosing evaluation metrics. It underscores the inadequacy of relying solely on traditional metrics like ROUGE for assessing the complex qualities desired in abstractive social media chat summaries, particularly faithfulness and semantic adequacy.

The state of evaluation in abstractive dialogue summarization reveals a significant gap. While BERTScore offers improvements over ROUGE by incorporating semantics [63], and LLM-as-judge approaches show promise for assessing aspects like coherence and faithfulness without references [48], no automatic metric currently provides a fully satisfactory substitute for nuanced human judgment. Human evaluation remains the gold standard but is hampered by cost, scalability issues, and potential inconsistencies, especially when detecting subtle faithfulness errors in fluent LLM outputs.[50] This evaluation bottleneck hinders rapid progress, as reliably measuring improvements in the very qualities abstractive methods aim for (like faithfulness and coherence beyond simple lexical match) remains difficult. Developing more reliable, scalable, and accurate automatic evaluation methods, particularly for faithfulness in the context of informal dialogue, is a critical area for future research.

## VII. Open Problems and Future Research Directions

Despite significant advancements driven by PLMs and LLMs, the field of context-aware social media chat summarization faces numerous open challenges and offers fertile ground for future research. Synthesizing the limitations discussed throughout this report reveals several key areas requiring further investigation.

**Synthesizing Current Limitations**

- **Faithfulness and Hallucination:** This remains arguably the most critical challenge.[1] Abstractive models, especially LLMs, frequently generate summaries containing information not present in or contradicting the source chat.[47] The tendency towards "circumstantial inference" in dialogue summarization, where models generate plausible but unsupported statements based on context, exacerbates this problem.[51] Reliably detecting and mitigating these hallucinations is paramount for building trustworthy summarization systems.

- **Coherence and Consistency:** Generating summaries that maintain a logical flow and internal consistency, especially when synthesizing information from long, multi-topic, and multi-participant conversations typical of social media, remains difficult.[29]

- **Effective Long-Context Handling:** While context window sizes are increasing, ensuring models effectively *utilize* information across the entire history, rather than just the beginning and end (mitigating the "lost-in-the-middle" effect [47]), is crucial for context-aware summarization of extended chats.[4]

- **Robust Evaluation:** The lack of reliable, scalable automatic evaluation metrics that correlate well with human judgments of summary quality (particularly faithfulness, coherence, and contextual relevance) hinders progress.[1]

- **Seamless Context Integration:** Developing effective mechanisms to fuse diverse contextual information – dialogue history, user profiles/preferences, and external knowledge sources – into the summarization process is an ongoing challenge.[2]

- **Handling Social Media Nuances:** Models need greater robustness to the specific linguistic characteristics of social media, including pervasive informality, rapidly evolving slang, ambiguous emojis, code-switching, noise, and platform-specific communication styles.[7]

- **Ethical Considerations (Bias, Fairness, Privacy, Toxicity):** Ensuring that summarization systems do not perpetuate societal biases present in the training data, treat all participants fairly, protect user privacy, and handle toxic content appropriately are critical ethical challenges that require technical solutions.[27]

- **Controllability and Personalization:** Enhancing the ability to generate summaries that meet specific user needs, adhere to constraints (e.g., length, style, focus), or are personalized based on user context remains an active area of research.[2]

**Proposing Novel Research Avenues**

Building upon the user's initial thoughts and the identified limitations, several promising and novel research directions emerge within context-aware social media chat summarization using BERT/LLMs:

(a) Improving Long-Context Handling for Social Media Dynamics:

Current long-context solutions often treat the input as a monolithic block or use generic segmentation. Research is needed on methods specifically tailored to the dynamic structure of social media conversations. This could involve:

● Developing models that explicitly recognize and adapt to frequent topic shifts, perhaps by dynamically segmenting the conversation based on semantic coherence.

● Investigating graph-based representations that model the conversational flow, participant interactions, and reply structures over long periods, potentially informing attention mechanisms.

● Designing hierarchical models that process information at utterance, turn, segment, and conversation levels, specifically tuned for chat characteristics.

● Refining retrieval mechanisms for retrieve-then-summarize pipelines to better suit chat dynamics, perhaps prioritizing recency while using semantic search to pull in crucial, topically relevant historical context, even if distant.

● Directly tackling the "lost-in-the-middle" problem [47] within summarization architectures, potentially through modified attention patterns or positional encoding schemes that encourage more uniform context utilization for chat data.

(b) Developing Methods Robust to Noise, Informality, and Code-Switching:

Enhancing model resilience to the non-standard and noisy nature of social media language is critical. Potential avenues include:

● Large-scale pre-training or continued pre-training on diverse, noisy social media corpora (e.g., leveraging datasets like WildChat [57] or large Reddit dumps [6]) to expose models to realistic language variations.

● Designing adaptive tokenization strategies or dedicated text normalization layers (potentially using style transfer techniques [56]) as pre-processing steps to handle slang, misspellings, and abbreviations before summarization.

● Exploring few-shot or zero-shot adaptation techniques that allow models to quickly learn and interpret new slang terms or emoji usage patterns specific to a particular community or platform.

● For code-switching, focusing on developing truly code-switch-aware multilingual models or pipelines that robustly identify language segments and apply appropriate linguistic processing, potentially integrating specialized lexicons or translation components.

(c) Creating Better Evaluation Metrics for Contextual Chat Summarization:

Addressing the evaluation bottleneck requires moving beyond traditional metrics. Research could focus on:

● Designing context-aware faithfulness metrics: These metrics would verify summary statements not just against the immediate source turns but against the relevant dialogue history required for interpretation, potentially identifying contradictions or unsupported claims that depend on earlier context.

● Developing coherence metrics based on conversational structure: Evaluating summary coherence not just based on linguistic fluency but also on how well it reflects the logical flow, topic shifts, and participant interactions of the original chat, possibly using graph-based analysis.

● Refining LLM-as-judge protocols for chat: Creating sophisticated prompting strategies or multi-agent frameworks [48] specifically designed to evaluate chat summaries, instructing the LLM evaluator to consider chat-specific nuances like informality, turn-taking, and implicit meaning when assessing quality aspects like faithfulness, coherence, and relevance.

● Metrics for robustness: Developing metrics that specifically quantify a summarization model's robustness to noise, information sparsity, or the presence of informal language elements in the source chat.

(d) Exploring Multimodal Social Media Chat Summarization:

Social media is increasingly multimodal, with images, GIFs, videos, and complex stickers often playing a crucial role in communication alongside text. Research is needed to:

● Develop multimodal fusion techniques capable of integrating visual and textual information effectively within the context of a chat conversation. How does an image or GIF modify the meaning or salience of adjacent text utterances?

● Investigate how non-textual elements contribute to the overall summary-worthy content of a chat segment.

● Create benchmark datasets for multimodal chat summarization, containing conversations with integrated visual elements and corresponding multimodal summaries.

● Explore the interplay between visual context and the interpretation of textual elements like emojis. [7]

(e) Investigating Personalization and User-Specific Summarization:

Tailoring summaries to individual users or specific perspectives offers significant value but requires further research:

● Developing models that can reliably infer user interests, information needs, or perspectives from their interaction history, explicit profiles, or conversational roles within the chat.[2]

● Advancing query-based summarization techniques [2] to handle complex, natural language queries reflecting diverse user goals in the context of informal chat.

● Exploring methods for generating viewpoint-specific summaries (e.g., summarizing a debate from one participant's perspective) or summaries that highlight information particularly relevant to a specific user based on their prior knowledge or role.

● Crucially, investigating and implementing robust privacy-preserving techniques (see point f) is integral to any personalization effort involving potentially sensitive user data.

(f) Addressing Fairness, Bias, Privacy, and Toxicity:

Integrating ethical considerations directly into the model design and evaluation process is essential for responsible deployment on social media. Key research directions include:

● Developing techniques for auditing and mitigating demographic biases in summarization models, ensuring that summaries fairly represent contributions from diverse participants and do not amplify harmful stereotypes potentially learned from biased training data.[27] This might involve fairness-aware training objectives or post-processing steps.

● Designing privacy-preserving summarization algorithms. This could involve techniques like federated learning (training models on decentralized user data), differential privacy (adding noise to obscure individual contributions), generating summaries that inherently omit or abstract sensitive personal information, or exploring on-device summarization.

● Integrating toxicity detection [59] and mitigation strategies into the summarization pipeline. This could involve filtering toxic content before summarization, training models to summarize non-toxically even when toxicity is present in the source, or explicitly flagging potentially problematic content in the summary output.

● Researching the inherent trade-offs between summary informativeness, conciseness, and ethical constraints like fairness and privacy.

The interconnectedness of these open problems is apparent. For instance, advances in long-context handling [47] are necessary not only for better context integration but also for enabling more accurate context-aware evaluation.[48] Similarly, effectively handling linguistic noise and informality [7] is closely linked to fairness concerns [27], as non-standard language use might correlate with specific demographic groups whose contributions could be inadvertently down-weighted by models struggling with such language. Addressing these challenges often requires considering their interplay, pushing towards more holistic solutions.

The specific type of hallucination identified as "Circumstantial Inference" [51] poses a particularly challenging problem for social media chat summarization. The inherent ambiguity, implicitness, and reliance on shared context in informal chat provide fertile ground for LLMs to generate plausible-sounding statements that lack direct evidential support in the text. Models optimized for coherence might readily "fill in the gaps" in these ambiguous conversations, potentially leading to a high prevalence of this subtle yet damaging form of hallucination. Tackling this requires developing models with a better capacity for distinguishing grounded facts from plausible inferences and evaluation methods specifically attuned to detecting such errors.

Furthermore, the drive towards personalization and user-specific summaries [2] creates a direct tension with the critical need for privacy protection in the context of often personal and sensitive social media data. Achieving utility through personalized summaries while rigorously safeguarding user privacy necessitates careful navigation. This highlights the need for research that explicitly integrates privacy-enhancing technologies (like differential privacy, federated learning, secure multi-party computation, or careful anonymization strategies) into the design and training of personalized summarization models. Balancing these competing demands will be a defining challenge for the field moving forward.

## VIII. Conclusion

This report has surveyed the landscape of context-aware summarization for social media chat, focusing on the role of BERT and advanced NLP techniques, particularly Large Language Models. The field has rapidly evolved from extractive methods towards abstractive approaches, driven by the power of sequence-to-sequence models, Transformers, and now LLMs, which offer unprecedented fluency and zero/few-shot capabilities.[29] However, this shift brings the challenge of ensuring factual consistency (faithfulness) to the forefront.[48]

Summarizing social media chat presents unique and synergistic challenges stemming from its inherent informality (slang,

emojis, code-switching), noise, multi-party dynamics, frequent topic shifts, and the variable quality of user-generated content.[4] Addressing these requires models that are robust, adaptable, and deeply context-aware. Context-awareness itself is multi-faceted, potentially incorporating dialogue history, user information, and external knowledge, demanding sophisticated integration techniques.[2]

BERT and subsequent LLMs provide powerful tools, but their application is constrained by limitations in handling very long contexts effectively (specifically, the "lost-in-the-middle" problem [47]), their propensity for hallucination (especially plausible but unsupported "circumstantial inferences" in dialogue [51]), sensitivity to prompting [33], and the significant computational costs involved.[61] Furthermore, evaluating the quality of generated summaries remains a major bottleneck, as standard metrics like ROUGE correlate poorly with human judgments of abstractive quality, and reliably assessing faithfulness automatically is an unsolved problem.[45] Available datasets like SAMSum and DialogSum provide valuable resources, but may not fully capture the diversity of real-world social media platforms.[6]

Based on this analysis, several promising research directions emerge for the specific topic of Context-Aware Summarization of Social Media Chat Using BERT and Enhanced NLP Techniques:

1. Developing novel long-context handling mechanisms tailored to chat dynamics and mitigating known LLM utilization issues.
2. Creating models inherently robust to the noise, informality, and code-switching characteristic of social media language, potentially through targeted pre-training or adaptive techniques.
3. Designing more reliable automatic evaluation metrics focused on contextual relevance, coherence, and, critically, faithfulness, perhaps leveraging LLM-as-judge frameworks specifically adapted for chat nuances.
4. Exploring multimodal summarization to incorporate the rich visual information present in social media conversations.
5. Investigating personalized and user-specific summarization methods while rigorously addressing the associated privacy implications.
6. Explicitly integrating fairness, bias mitigation, and toxicity handling into summarization models to ensure responsible deployment in sensitive social media contexts.

Addressing these open problems holds significant potential for advancing the field, enabling the development of more accurate, reliable, and useful tools for navigating the vast and complex landscape of online conversations. The intersection of context-awareness, social media's unique challenges, and the capabilities (and limitations) of modern LLMs provides a rich and impactful area for continued research.

**Works cited**

1. CADS: A Systematic Literature Review on the Challenges of Abstractive Dialogue Summarization - arXiv, accessed April 29, 2025, https://arxiv.org/html/2406.07494v3
2. Instructive Dialogue Summarization with Query Aggregations - ACL Anthology, accessed April 29, 2025, https://aclanthology.org/2023.emnlp-main.474.pdf
3. Novel framework for dialogue summarization based on factual-statement fusion and dialogue segmentation - PMC - PubMed Central, accessed April 29, 2025, https://pmc.ncbi.nlm.nih.gov/articles/PMC11020369/
4. www.ijcai.org, accessed April 29, 2025, https://www.ijcai.org/proceedings/2022/0764.pdf
5. FREDSum: A Dialogue Summarization Corpus for French Political Debates - ACL Anthology, accessed April 29, 2025, https://aclanthology.org/2023.findings-emnlp.280.pdf
6. +64 Summarization Datasets - NLP Database - Metatext - AI, accessed April 29, 2025, https://metatext.io/datasets-list/summarization-task
7. Emojis Decoded: Leveraging ChatGPT for Enhanced Understanding in Social Media Communications - arXiv, accessed April 29, 2025, https://arxiv.org/html/2402.01681v3
8. Tailoring Generative AI Chatbots for Multiethnic Communities in Disaster Preparedness Communication: Extending the CASA Paradigm - arXiv, accessed April 29, 2025, https://arxiv.org/html/2406.08411v2
9. [2301.12303] Presence of informal language, such as emoticons, hashtags, and slang, impact the performance of sentiment analysis models on social media text? - arXiv, accessed April 29, 2025, https://arxiv.org/abs/2301.12303

10.	arxiv.org, accessed April 29, 2025, https://arxiv.org/pdf/2301.12303

11.	(PDF) A Systematic Literature Review on Social Media Slang Analytics in Contemporary Discourse - ResearchGate, accessed April 29, 2025, https://www.researchgate.net/publication/375735233_A_Systematic_Literature_Review_on_Social_Media_Slang_Analytics_in_Contemporary_Discourse

12.	arXiv:2504.08024v1 [cs.CL] 10 Apr 2025, accessed April 29, 2025, https://arxiv.org/pdf/2504.08024

13.	Exploring the Efficacy of Large Language Models in Summarizing Mental Health Counseling Sessions: Benchmark Study, accessed April 29, 2025, https://mental.jmir.org/2024/1/e57306

14.	An Exploratory Study on Long Dialogue Summarization: What Works and What's Next - Microsoft, accessed April 29, 2025, https://www.microsoft.com/en-us/research/wp-content/uploads/2021/09/4069_Paper-1.pdf

15.	Extractive vs Abstractive Summarization in Healthcare, accessed April 29, 2025, https://www.abstractivehealth.com/article/extractive-vs-abstractive-summarization-in-healthcare

16.	Exploring the Extractive Method of Text Summarization - Analytics Vidhya, accessed April 29, 2025, https://www.analyticsvidhya.com/blog/2023/03/exploring-the-extractive-method-of-text-summarization/

17.	Contextualizing Search Queries In-Context Learning for Conversational Rewriting with LLMs - arXiv, accessed April 29, 2025, https://arxiv.org/html/2502.15009v1

18.	A Context-aware Framework for Translation-mediated Conversations - arXiv, accessed April 29, 2025, https://arxiv.org/html/2412.04205v1

19.	History-Aware Conversational Dense Retrieval - arXiv, accessed April 29, 2025, https://arxiv.org/html/2401.16659v1

20.	Interpersonal Memory Matters: A New Task for Proactive Dialogue Utilizing Conversational History - arXiv, accessed April 29, 2025, https://arxiv.org/html/2503.05150v1

21.	Context-Aware LLM Translation System Using Conversation ..., accessed April 29, 2025, https://aclanthology.org/2024.wmt-1.102/

22.	CADGE: Context-Aware Dialogue Generation Enhanced with Graph-Structured Knowledge Aggregation - ACL Anthology, accessed April 29, 2025, https://aclanthology.org/2024.inlg-main.31.pdf

23.	Long Dialog Summarization: An Analysis - arXiv, accessed April 29, 2025, https://arxiv.org/html/2402.16986v1

24.	Cutting through social media noise: Strategies to navigating a cluttered social media landscape - Sprout Social, accessed April 29, 2025, https://sproutsocial.com/insights/social-media-noise/

25.	Short-form video, user-generated content are among top social media trends for 2025, accessed April 29, 2025, https://www.inma.org/blogs/social-media/post.cfm/short-form-video-user-generated-content-are-among-top-social-media-trends-for-2025

26.	Social Context Summarization using User-generated Content and Third-party Sources | Request PDF - ResearchGate, accessed April 29, 2025, https://www.researchgate.net/publication/322106647_Social_Context_Summarization_using_User-generated_Content_and_Third-party_Sources

27.	Summarizing User-generated Textual Content: Motivation and Methods for Fairness in Algorithmic Summaries - arXiv, accessed April 29, 2025, https://arxiv.org/pdf/1810.09147

28.	NUT-RC: Noisy User-generated Text-oriented Reading Comprehension - ACL Anthology, accessed April 29, 2025, https://aclanthology.org/2020.coling-main.242.pdf

29.	A Comprehensive Survey on Automatic Text Summarization with Exploration of LLM-Based Methods - arXiv, accessed April 29, 2025, https://arxiv.org/html/2403.02901v2

30.	Large-scale Summarization of Chat Transcripts in the Absence of Annotated Summaries - ACL Anthology, accessed April 29, 2025, https://aclanthology.org/2024.icnlsp-1.12.pdf

31.	A Systematic Survey of Text Summarization: From Statistical Methods to Large Language Models - arXiv, accessed April 29, 2025, https://arxiv.org/html/2406.11289v1

32.	arXiv:2406.15809v4 [cs.CL] 24 Jan 2025 - Garima Chhikara, accessed April 29, 2025, https://garimachhikara128.github.io/papers/LLM_Summarization.pdf

33.	Zero-shot Conversational Summarization Evaluations with small Large Language Models - arXiv, accessed April 29, 2025, https://arxiv.org/pdf/2311.18041

34.     Abstractive Text Summarization: State of the Art, Challenges, and Improvements - arXiv, accessed April 29, 2025, https://arxiv.org/html/2409.02413v1

35.     Summarization - Hugging Face, accessed April 29, 2025, https://huggingface.co/docs/transformers/tasks/summarization

36.     Extractive vs. Abstractive Summarization: How Does it Work? - Prodigal, accessed April 29, 2025, https://www.prodigaltech.com/blog/extractive-vs-abstractive-summarization-how-does-it-work

37.     On model selection for text summarization - Mozilla.ai, accessed April 29, 2025, https://blog.mozilla.ai/on-model-selection-for-text-summarization/

38.     Abstractive vs. Extractive Summarization: An Experimental Review - MDPI, accessed April 29, 2025, https://www.mdpi.com/2076-3417/13/13/7620

39.     Dialogue Summarization: A Deep Learning Approach - Analytics Vidhya, accessed April 29, 2025, https://www.analyticsvidhya.com/blog/2021/02/dialogue-summarization-a-deep-learning-approach/

40.     A Comprehensive Survey on Automatic Text Summarization with Exploration of LLM-Based Methods - arXiv, accessed April 29, 2025, https://arxiv.org/pdf/2403.02901

41.     Dialogue Summarization with Flan-T5 | summarization-trial-t5 – Weights & Biases - Wandb, accessed April 29, 2025, https://wandb.ai/events/summarization-trial-t5/reports/Dialogue-Summarization-with-Flan-T5--VmlldzozNjg5NTU5

42.     Automatic Short Text Summarization Techniques in Social Media Platforms - MDPI, accessed April 29, 2025, https://www.mdpi.com/1999-5903/15/9/311

43.     [2410.15962] Systematic Exploration of Dialogue Summarization Approaches for Reproducibility, Comparative Assessment, and Methodological Innovations for Advancing Natural Language Processing in Abstractive Summarization - arXiv, accessed April 29, 2025, https://arxiv.org/abs/2410.15962

44.     In-context Learning of Large Language Models for Controlled Dialogue Summarization: A Holistic Benchmark and Empirical Analysis - ACL Anthology, accessed April 29, 2025, https://aclanthology.org/2023.newsum-1.6.pdf

45.     Daily Papers - Hugging Face, accessed April 29, 2025, https://huggingface.co/papers?q=SAMSum%20Corpus

46.     Text Summarization | Papers With Code, accessed April 29, 2025, https://paperswithcode.com/task/text-summarization/codeless?page=13&q=

47.     On Context Utilization in Summarization with Large Language Models - arXiv, accessed April 29, 2025, https://arxiv.org/html/2310.10570v3

48.     Faithful, Unfaithful or Ambiguous? Multi-Agent Debate with Initial Stance for Summary Evaluation - arXiv, accessed April 29, 2025, https://arxiv.org/html/2502.08514v2

49.     On Assessing the Faithfulness of LLM-generated Feedback on Student Assignments, accessed April 29, 2025, https://educationaldatamining.org/edm2024/proceedings/2024.EDM-short-papers.49/index.html

50.     STORYSUMM: Evaluating Faithfulness in Story Summarization - ACL Anthology, accessed April 29, 2025, https://aclanthology.org/2024.emnlp-main.557.pdf

51.     Analyzing LLM Behavior in Dialogue Summarization: Unveiling Circumstantial Hallucination Trends - arXiv, accessed April 29, 2025, https://arxiv.org/html/2406.03487v1

52.     Analyzing LLM Behavior in Dialogue Summarization: Unveiling Circumstantial Hallucination Trends - ACL Anthology, accessed April 29, 2025, https://aclanthology.org/2024.acl-long.677.pdf

53.     Building high-quality datasets for abstractive text summarization - DiVA portal, accessed April 29, 2025, https://www.diva-portal.org/smash/get/diva2:1563580/FULLTEXT01.pdf

54.     Large Language Models: A Survey - arXiv, accessed April 29, 2025, https://arxiv.org/html/2402.06196v2

55.     Mixture of In-Context Experts Enhance LLMs' Long Context Awareness - NIPS papers, accessed April 29, 2025, https://proceedings.neurips.cc/paper_files/paper/2024/file/91315fbb83ce353ae5538cba395f70d1-Paper-Conference.pdf

56.     A Review of Text Style Transfer using Deep Learning - arXiv, accessed April 29, 2025, https://arxiv.org/pdf/2109.15144

57.     Open data - Ai2, accessed April 29, 2025, https://allenai.org/open-data

58.     lmsys/lmsys-chat-1m · Datasets at Hugging Face, accessed April 29, 2025, https://huggingface.co/datasets/lmsys/lmsys-chat-1m

59. Combating Toxic Language: A Review of LLM-Based Strategies for Software Engineering, accessed April 29, 2025, https://arxiv.org/html/2504.15439v1

60. Efficacy of Context Summarization Techniques on Large Language Model Chatbots - DiVA portal, accessed April 29, 2025, http://www.diva-portal.org/smash/get/diva2:1886192/FULLTEXT01.pdf

61. QMSum: A New Benchmark for Query-based Multi-domain Meeting Summarization | Request PDF - ResearchGate, accessed April 29, 2025, https://www.researchgate.net/publication/352365008_QMSum_A_New_Benchmark_for_Query-based_Multi-domain_Meeting_Summarization

62. Intent-Aware Dialogue Generation and Multi-Task Contrastive Learning for Multi-Turn Intent Classification - arXiv, accessed April 29, 2025, https://arxiv.org/html/2411.14252v1

63. Improving Neural Abstractive Summarization via Reinforcement Learning with BERTScore - CS229 - Stanford University, accessed April 29, 2025, https://cs229.stanford.edu/proj2019aut/data/assignment_308832_raw/26632588.pdf

64. A Primer on Large Language Models and their Limitations - arXiv, accessed April 29, 2025, https://arxiv.org/html/2412.04503v1

65. On Assessing the Faithfulness of LLM-generated Feedback on Student Assignments - Educational Data Mining, accessed April 29, 2025, https://educationaldatamining.org/edm2024/proceedings/2024.EDM-short-papers.49/2024.EDM-short-papers.49.pdf

66. How to evaluate a summarization task - OpenAI Cookbook, accessed April 29, 2025, https://cookbook.openai.com/examples/evaluation/how_to_eval_abstractive_summarization

67. README.md - THUDM/LongBench · GitHub, accessed April 29, 2025, https://github.com/THUDM/LongBench/blob/main/LongBench/README.md

68. Dialog Summarization - Kaggle, accessed April 29, 2025, https://www.kaggle.com/datasets/marawanxmamdouh/dialogsum

69. Public Data Sources - Text - Madrona Venture Labs, accessed April 29, 2025, https://www.madronavl.com/launchable/public-data-sources-text

70. Evaluating Small Language Models for News Summarization: Implications and Factors Influencing Performance - arXiv, accessed April 29, 2025, https://arxiv.org/html/2502.00641v2

71. Evaluating LLMs and Pre-trained Models for Text Summarization Across Diverse Datasets - arXiv, accessed April 29, 2025, https://arxiv.org/pdf/2502.19339

72. LLM evaluation metrics and methods - Evidently AI, accessed April 29, 2025, https://www.evidentlyai.com/llm-guide/llm-evaluation-metrics

73. Balancing Lexical and Semantic Quality in Abstractive Summarization - ACL Anthology, accessed April 29, 2025, https://aclanthology.org/2023.acl-short.56.pdf

74. The official repo for paper, LLMs-as-Judges: A Comprehensive Survey on LLM-based Evaluation Methods. - GitHub, accessed April 29, 2025, https://github.com/CSHaitao/Awesome-LLMs-as-Judges

75. Learning to Verify Summary Facts with Fine-Grained LLM Feedback - arXiv, accessed April 29, 2025, https://arxiv.org/html/2412.10689v1

76. Hallucinations in LLMs and Resolving Them: A Holistic Approach - SciTePress, accessed April 29, 2025, https://www.scitepress.org/Papers/2025/130945/130945.pdf