# *ChatGPT in Research: Variable Extraction and Researcher Protection in the Context of Child Sexual Abuse*

Noreen Naranjos Velazquez, IU International University of Applied Sciences, Germany

**Abstract**

In the context of research on childhood sexual abuse (CSA), researchers face significant challenges due to the emotionally distressing nature of sensitive data (Williamson et al., 2020). The use of ChatGPT (OpenAI, 2023) offers valuable support in this area. This artificial intelligence (AI) tool facilitates efficient data processing while maintaining emotional distance by converting qualitative content into quantifiable formats. This approach not only aids in statistical analysis but also reduces the emotional burden on researchers (van Manen, 2023). The inter-coder reliability of this method has been evaluated and largely confirmed in various forms (Naranjos Velazquez, 2024; in press). In a comparative study, the results of AI models ChatGPT 3.5, ChatGPT 4 and ChatGPT 4o were analysed. Increased consistency was observed beginning with the ChatGPT 4 model, further highlighting the reliability of ChatGPT in processing sensitive information. This presentation explores the ethical and practical implications of AI use in research and discusses the limitations of this AI tool (Naranjos Velazquez, 2023[a]; in press).

Keywords: Artificial Intelligence, AI, ChatGPT, CSA, Inter-coder Reliability, Secondary Traumatization, Sexual Abuse of Children

iafor

The International Academic Forum
www.iafor.org

**Introduction**

This study explores methodological and emotional challenges in researching childhood sexual abuse (CSA) and proposes innovative AI-assisted approaches. To address these challenges, this study proposes the integration of network analysis and the use of artificial intelligence tools like ChatGPT as innovative approaches to process qualitative data more efficiently and objectively. Building on these challenges, the article introduces the integration of network analysis as a novel approach to analyse survivors' self-reports of CSA. In addition, the second part of the paper explores the application of ChatGPT, an artificial intelligence (AI) tool, for extracting network data from mentioned narratives and statistically comparing results between human coder and AI. The paper demonstrates how this AI tool can support researchers by converting qualitative data into quantifiable data. The discussion also critically addresses ethical dimensions and practical implications of employing supportive AI in CSA research.

**Challenges and Approaches in CSA Research**

*Key Challenges*

In contrast to official statistics, self-reports of CSA reveal abuse rates up to 30 times higher than official reports. This self-reports are very important, because CSA is particularly challenging to detect, and it often goes unreported until adulthood, with estimates showing that 60-80% of survivors disclose their abuse only later in life (Naranjos Velazquez, in press[a;b]). The delayed disclosure of CSA often precludes criminal justice action, as statutes of limitations may void legal consequences. Consequently, early intervention and therapeutic support are crucial yet rarely accessible in time (Naranjos Velazquez, in press[a]). One of the second main challenges in this field is the emotional impact on researchers who engage with detailed survivor accounts. This work can lead to secondary trauma, where researchers experience distress like that of survivors (Williamson et al., 2020).

*Theoretical Framework: A Socio-Ecological Perspective*

The disclosure of CSA involves complex social dynamics. Using a social-ecological model, in this study the roles of individual, familial, and societal factors in disclosure processes were focused using ChatGPT (OpenAI, 2023), a tool of artificial intelligence (AI) on self-reports of adult survivors of CSA (Naranjos Velazquez, in press[a]). ChatGPT offers a way to manage these distressing narratives, allowing researchers to process qualitative data more objectively while reducing emotional strain (Naranjos Velazquez, in press[a]; van Manen, 2023).

**Research Design: Network Analysis and AI Integration**

*Data Sources: Self-Reports of CSA*

The source of the data was self-reports collected from the "platform of histories" of the Independent Inquiry into Child Sexual Abuse (UKASK, 2024), with a sample size ($N = 113$) that includes 77 female and 36 male survivors. The UKASK platform serves as a repository for documenting the personal narratives of survivors of child sexual abuse. It aims to raise public awareness, provide insights into the survivors' experiences, and support systemic reforms by sharing these powerful stories. The platform "Geschichten, die zählen" (Stories that matter) was created to serve as a place of remembrance and to honor the life

achievements of those affected by sexual abuse, ensuring that their experiences are not forgotten. These German self-reports are openly accessible on the UKASK (2024) website. They were published based on submissions voluntarily provided by survivors. The reports were documented by UKASK through various means, including confidential hearings and written submissions. All reports were anonymized to protect the identities of the survivors, and their publication was approved with explicit consent. These self-reports provide a valuable qualitative dataset for analyzing relational dynamics and contextual variables in CSA narratives (Naranjos Velazquez, 2024; 2023[b]).

## Network Analysis of Survivor Social Circles

In network analysis, particularly in egocentric network analysis the social circles surrounding survivors can be studied (Naranjos Velazquez, in press[b]). This approach categorizes connections by strength and network size (Perry et al., 2018, p. 160); revealing essential patterns such as relationships with perpetrators and silent witnesses. For example, based on the conceptual framework of a social-ecological model, different levels of relational strength where the connection to family members may be strong but strained, while weaker ties of strangers might still play significant roles (Naranjos Velazquez, in press[a]). From the survivor's perspective, the size of the network is calculated as the total number of identified perpetrators or bystander (Naranjos Velazquez, 2023[b], pp. 91-94, in press[b]).

## Application of Personal Network Analysis

To illustrate the relevance of network analysis in CSA research, initial findings based on manually coded self-reports are presented. These analyses emphasize the importance of understanding relational dynamics in the social networks surrounding survivors' members (Naranjos Velazquez, in press[b]). In terms of gender and silent mothers, statistical tools such as the Chi-square test ($\chi^2$ = 9.12, $\phi$ = -.28, $df$ = 1, $p$ < .01) revealed significant gender differences in relational dynamics with silent witnesses, particularly mothers (Field, 2018, pp. 838-839). These findings demonstrate that gender significantly influences interaction patterns with key network members (Naranjos Velazquez, in press[b]).

Regarding the strength of ties, Fisher's Exact Test ($p$ < .001, $\phi$ = .50) was applied to assess the strength of ties between perpetrators and survivors, as well as the presence of silent mothers (Field, 2018, p. 839). These results highlight significant differences in relational strength depending on whether the mother assumes a silent or non-silent role, reflecting the intricate network structures surrounding survivors (Naranjos Velazquez, in press[b]).

Finally, with respect to network size, Spearman's Rho ($\rho$ = 0.19, $p$ = .04, 95% $CI$ [.00, .37]) identified a modest but statistically significant correlation between the network size of perpetrators and presence of silent individuals which, according to Cohen's guidelines, represents a small effect size (Cohen, 1988; Field, 2018, p. 344). These findings suggest that as the number of perpetrators increases, there is a slight tendency for the number of silent individuals in the network to grow (Naranjos Velazquez, in press[b]).

**Methodology and Data Analysis: Using AI for Variable Extraction**

*Using ChatGPT for Data Extraction*

ChatGPT's natural language processing (NLP) capabilities were employed to extract critical variables from survivor narratives efficiently (OpenAI, 2023). Using OpenAI´s Playground API (OpenAI, 2024), the data was processed on a secure server hosted by IU International University of Applied Science. This approach was necessary due to OpenAI's guidelines and the highly sensitive nature of the self-reports on CSA used in this study, which made it impossible to conduct the analysis using the publicly available version of ChatGPT. The publicly available version of ChatGPT blocks the processing of such prompts to ensure compliance with ethical and safety standards (OpenAI, 2023). Additionally, the API version allows for the specification of roles such as 'user,' 'system,' and 'assistant,' enabling more controlled and context-specific interactions during the data extraction process. To further structure the analysis, role-based prompt engineering was employed.

*Role-Based Prompt Engineering for Sensitive Narratives*

In this study, prompts were designed to utilize the OpenAI Playground API's ability to define specific roles such as "user," "system," and "assistant" (OpenAI, 2024). The system role was used to define the tone and scope of the analysis, instructing ChatGPT to process sensitive narratives on childhood sexual abuse with scientific focus and empathy. An example of a system prompt provided instructions like: *You are a scientist specializing in the analysis of reports on childhood sexual abuse. Your task is to identify key contexts while considering the effects of these experiences on the victims.*

As documented in OpenAI's guidelines for prompt engineering (OpenAI, 2024), the user role was designed to present specific questions for variable extraction, such as identifying the perpetrator, the presence of silent witnesses, or the context in which the abuse occurred. This approach provided clear instructions for each case. An example of a user prompt included questions like:
  * *Who is the perpetrator?*
  * *Were other forms of abuse (physical or psychological) also mentioned?*
  * *Name individuals who knew about the abuse but remained silent, specifying the context (e.g., family, school, religious community).*

The assistant role, representing ChatGPT's responses, structured the extracted data by listing variables or noting missing information. For instance, in one case study, ChatGPT responded with:
  * *Yes, physical abuse was mentioned.*
  * *The perpetrator was a teacher.*

This interaction model within the API enabled structured and context-specific responses, ensuring alignment with the study's requirements. To evaluate the effectiveness of these prompts, both zero-shot and few-shot learning approaches were compared (OpenAI, 2024).

*Comparison of Zero-Shot and Few-Shot Learning for Variable Extraction*

Specific prompts were designed to identify abuse contexts, silent witnesses, and active respondents. The analysis utilized both zero-shot learning and few-shot learning to assess

ChatGPT's performance in variable extraction (OpenAI, 2024). In the zero-shot learning approach, ChatGPT was provided with no prior examples and relied solely on the prompts to extract variables. In contrast, the few-shot learning approach included five examples, selected by ChatGPT itself from the dataset coded by humans (Naranjos Velazquez, in press[b]), based on the guideline to exemplify the widest possible range of variable expressions for the present study. For instance, the variable 'perpetrator' displayed diverse manifestations, as did other variables. The prompt explicitly instructed ChatGPT to compile the five examples to represent the broadest possible spectrum of these expressions, ensuring that the selected examples captured a wide variety of characteristics across the dataset. According to OpenAI's best practices for prompt engineering extraction (OpenAI, 2024), the use of carefully selected examples in few-shot learning can enhance the model's ability to generalize and adapt to complex datasets, thereby potentially improving the consistency of variable.

### Key Variables Extracted for CSA Research

Through the carefully designed prompts, ChatGPT extracted key variables, including "perpetrator identity", "silent witnesses", and "active respondents". Additional variables, such as the "context of violence", "presence of silent witnesses", "active individuals who knew about the abuse and took action", and the "age of survivors during the abuse" were also included, based on Naranjos Velazquez (in press[b]). These variables were critical for understanding the dynamics of abuse contexts and their broader implications for CSA research. By leveraging OpenAI's Playground API (OpenAI, 2024), the model ensured a consistent and systematic approach to variable extraction.

### Agreement Analysis Between Human Coder and ChatGPT

To evaluate the agreement between human coder and ChatGPT in extracting critical variables, inter-coder reliability was assessed using Cohen's Kappa ($\kappa$), as described by Gwet (2008). Both zero-shot learning (no prior examples provided) and few-shot learning (five examples provided) approaches were applied to compare consistency in variable extraction. According to Cohen (1960), Kappa values are interpreted as follows: values below 0.20 are considered poor, 0.21 to 0.40 fair, 0.41 to 0.60 moderate, 0.61 to 0.80 good, 0.81 to 0.99 very good, and 1.0 indicates perfect agreement. This methodological design allowed for an evaluation of ChatGPT's ability to generalize and adapt to complex datasets. All statistical analyses were conducted using SPSS Statistics (IBM Corp., 2021), with statistical significance set at $\alpha = 0.05$.

## Results

This study organizes its findings around the performance of zero-shot and few-shot learning techniques, evaluating ChatGPT's reliability and adaptability in analyzing sensitive CSA narratives (UKASK, 2024).

### Zero-Shot Learning Performance

In the zero-shot learning approach, ChatGPT demonstrated varied levels of inter-coder reliability across different variables. The agreement for identifying "perpetrator identity" was very high (Cohen's $\kappa = 0.88$, $p < .001$ for GPT 3.5; $\kappa = 0.94$, $p < .001$ for GPT 4 and GPT 4.o) and consistent across all tested versions. The agreement for "silent person" improved incrementally with more advanced versions, from $\kappa = 0.26$ ($p < .01$) for GPT 3.5 to $\kappa = 0.30$

(*p* < .001) for GPT 4 and reaching κ = 0.64 (*p* < .001) for GPT 4.o, indicating moderate agreement with the latest model. However, performance on "active person" exhibited a decline, with Cohen's κ dropping from 0.53 for GPT 3.5 to 0.28 for GPT 4, and further to 0.26 for GPT 4.o, reflecting lower consistency (p < .001) in identifying this variable. For "context of violence," agreement was moderate to good, improving from κ = 0.44 (*p* < .01) for GPT 3.5 to κ = 0.73 (*p* < .001) for GPT 4, and peaking at κ = 0.75 (*p* < .001) for GPT 4.o. The variable "age during CSA" showed fair agreement across all models, with κ values ranging from 0.29 (*p* = .01) for GPT 3.5, 0.28 (*p* < .001) for GPT 4, to 0.26 (*p* < .001) for GPT 4.o.

Table 1: Agreement Between Human Coder and ChatGPT (Zero-Shot Learning)

| Variable | GPT 3.5 | | | GPT 4 | | | | GPT 4.o | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Cohens Kappa[a] (κ) | *p*-value | *SE* | Cohens (κ) | Kappa[a] | *p*-value | *SE* | Cohens Kappa[a] (κ) | *p*-value | *SE* |
| Context of violence | 0.44 | < .01 | .14 | 0.73 | | < .001 | .10 | 0.75 | < .001 | .09 |
| Perpetrator | 0.88 | < .001 | .05 | 0.94 | | < .001 | .04 | 0.94 | < .001 | .04 |
| Silent person | 0.26 | < .01 | .10 | 0.30 | | < .001 | .09 | 0.64 | < .001 | .09 |
| Active person | 0.53 | < .001 | .16 | 0.28 | | < .001 | .10 | 0.26 | < .001 | .10 |
| Age during CSA | 0.29 | .010 | .16 | 0.28 | | < .001 | .10 | 0.26 | < .001 | .10 |

*Note.* κ = Cohen's Kappa value (α = .05), a measure of inter-coder agreement Cohen, 1960); *SE* = Standard Error. [a]According to Cohen (1960), values below 0.20 are considered poor, 0.21 to 0.40 fair, 0.41 to 0.60 moderate, 0.61 to 0.80 good, 0.81 to 0.99 very good, and 1.0 perfect agreement.

### *Few-Shot Learning Performance*

In the few-shot learning approach, ChatGPT's inter-coder reliability varied depending on the variable and version. For "perpetrator identity" a very high agreement was observed across all versions, peaking with GPT 4.o (Cohen's κ = .93, p < .001). The agreement for "silent person" slightly declined from GPT 3.5 (κ = .45, p < .001) to GPT 4.o (κ = .40, p < .001) and remained relatively low. Performance on "active person" was consistently poor across (*p* < .05) all versions, with κ values ranging from .10 (GPT 4) to .16 (GPT 3.5). The "context of violence" variable showed fluctuating levels of agreement. It demonstrated moderate agreement with GPT 3.5 (κ = .45, *p* < .01) and reached its highest level with GPT 4.o (κ = .58, *p* < .001). However, agreement dropped substantially for GPT 4 (κ = .15, no significant p-value), highlighting inconsistencies across model versions. The variable "age during CSA" demonstrated fair agreement, with κ values ranging from .20 (GPT 4) to .41 (GPT 3.5, p < .05).

Table 2: Agreement Between Human Coder and ChatGPT (Few-Shot Learning)

| Variable | GPT 3.5 | | | GPT 4 | | | | GPT 4.o | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Cohens Kappa[a] (κ) | *p*-value | *SE* | Cohens (κ) | Kappa[a] | *p*-value | *SE* | Cohens Kappa[a] (κ) | *p*-value | *SE* |
| context of violence | 0.45 | < .01 | .15 | 0.15 | | .12 | .08 | 0.58 | < .001 | .12 |
| perpetrator | 0.80 | < .001 | .07 | 0.80 | | < .001 | .08 | 0.93 | < .001 | .05 |
| silent person | 0.45 | < .001 | .10 | 0.13 | | .13 | .10 | 0.40 | < .001 | .11 |
| active person | 0.16 | .01 | .10 | 0.10 | | < .01 | .05 | 0.13 | < .001 | .06 |
| age during CSA | 0.41 | < .01 | .17 | 0.20 | | .08 | .15 | 0.37 | < .01 | .16 |

*Note.* κ = Cohen's Kappa value (α = .05), a measure of inter-coder agreement Cohen, 1960); *SE* = Standard Error. [a]According to Cohen (1960), values below 0.20 are considered poor, 0.21 to 0.40 fair, 0.41 to 0.60 moderate, 0.61 to 0.80 good, 0.81 to 0.99 very good, and 1.0 perfect agreement.

**Discussion**

*Comparison of Zero-Shot and Few-Shot Learning*

The analysis highlights notable differences in the performance of zero-shot and few-shot learning approaches when applied to variable extraction from CSA survivor narratives. Zero-shot learning demonstrated consistent reliability for straightforward variables like "perpetrator identity," showcasing its ability to handle clearly defined and less context-dependent information. However, its performance decreased significantly for nuanced variables such as "silent person" and "active person," which require an understanding of complex relational dynamics embedded in survivor narratives. These findings align with prior research emphasizing the challenges AI faces in addressing subtle and context-specific information in sensitive contexts (Naranjos Velazquez, in press[a;b]; Williamson et al., 2020).

In contrast, few-shot learning incorporated five carefully selected contextual examples to improve performance. While OpenAI's best practices for few-shot learning recommend using diverse examples to maximize adaptability (OpenAI, 2024), the limited number of examples in this analysis constrained the model's ability to generalize to more complex variables. "Perpetrator identity" continued to show consistently high agreement across versions, but nuanced variables like "silent person" and "active person" remained problematic. For instance, providing contextual examples only minimally enhanced performance, suggesting that few-shot learning's benefits are highly variable-specific and less effective for extracting complex relational dynamics (Naranjos Velazquez, in press[b]). The findings suggest that neither approach is universally sufficient for addressing the full spectrum of variables in CSA research. Zero-shot learning performs well for straightforward variables, while few-shot learning shows limited improvements for more nuanced variables, highlighting the need for alternative or hybrid approaches. Additionally, the discrepancies in performance across model versions underscore the importance of iterative refinement in AI methodologies to better address the challenges of sensitive and complex narratives (OpenAI, 2023; van Manen, 2023).

*Methodological Limitations and Implications*

Despite the advantages of AI-assisted methods like ChatGPT in analysing sensitive data in CSA research, several methodological limitations must be considered. First, the findings indicate that both zero-shot and few-shot learning approaches struggle to extract nuanced variables such as "silent person" and "active person," even when contextual examples are provided (Naranjos Velazquez, in press[a;b]). This limitation suggests that the models face difficulties abstracting complex social and interpersonal dynamics often implicitly embedded in survivor narratives. The lower Cohen's kappa values for these variables (e.g., $\kappa = .10$–$.16$ for "active person") underscore this issue, aligning with prior studies that highlight challenges in modelling relational variables in sensitive contexts). Second, the use of predefined prompts poses a risk of unintended bias due to the selection of specific examples, potentially limiting the generalizability of results (OpenAI, 2024). Particularly in few-shot learning, there is a risk that the examples may not capture the full variability of the dataset, leading to systematic over- or underestimation of certain variables (OpenAI, 2023; van Manen, 2023). Moreover, discrepancies observed between models (e.g., GPT 3.5, GPT 4, and GPT 4.o) highlight inconsistencies in AI-assisted analyses, even within the same technological framework (Naranjos Velazquez, 2024). Additionally, the dataset, composed of 113 self-reports available on the UKASK platform (UKASK, 2024), while valuable, poses limitations in

generalizability due to its reliance on survivor narratives and the specific context in which they were collected. Although these reports provide a rich source of qualitative data, their scope is inherently limited to the individuals who chose to disclose their experiences. Moreover, the sensitive nature of the data necessitated strict adherence to ethical standards. To this end, all analyses were conducted on a secure internal server hosted by IU International University of Applied Sciences, leveraging OpenAI's API for structured and controlled processing (OpenAI, 2023). These findings underline the need for methodological innovations to address the limitations of both zero-shot and few-shot learning. Future research should explore the integration of hybrid approaches, combining AI's computational efficiency with human expertise. For example, expanding the number and diversity of examples in few-shot learning could enhance the model's ability to generalize. Similarly, incorporating domain-specific training data could improve AI's performance for complex relational variables, such as "silent witnesses" and "active individuals" (Naranjos Velazquez, in press[a;b]). Future advancements in AI methodologies must prioritize expanding the adaptability of few-shot learning while ensuring that ethical considerations remain central to the development of tools for analyzing sensitive data. By addressing these challenges, researchers can enhance both the effectiveness and the ethical robustness of AI-assisted approaches in CSA research.

**Conclusion**

Using a socio-ecological model helps capture the multifaceted influences surrounding childhood sexual abuse, from individual to societal levels. Analysing personal networks reveals the roles of both perpetrators and silent witnesses, offering deeper insights into abuse dynamics. Survivor self-reports provide essential variables, ensuring data directly reflects their experiences. The integration of ChatGPT's natural language processing capabilities has proven valuable in extracting key variables, enabling more systematic analyses while maintaining emotional safety for researchers. Strong inter-coder agreement on specific variables, especially with provided examples, enhances data reliability. However, challenges in extracting nuanced variables, such as silent witnesses or active bystanders, highlight the limitations of AI in addressing complex relational dynamics. Strict ethical standards and data protection are upheld to ensure survivors' privacy and dignity. Future research should explore hybrid approaches that combine AI efficiency with human expertise to address these methodological challenges.

# References

Cohen, J. (1988). Statistical power analysis for the behavioral sciences. (2nd ed.). Hoboken: Taylor and Francis.

Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20, 37–46.

Field, A. (2018). Discovering Statistics Using IBM SPSS Statistics. (5th ed.). London: SAGE Publications.

Gwet, K. L. (2008). Computing inter-rater reliability and its variance in the presence of high agreement. *The British Journal of Mathematical and Statistical Psychology, 61*(1), 29–48. https://doi.org/10.1348/000711006X126600

IBM Corp. (2021). IBM SPSS Statistics for Windows (Version 28.0) [Computer software]. New York: IBM Corp.

Naranjos Velazquez, N. (2023a). Use of ChatGPT in networks of early childhood interventions. Rostock: Universität Rostock. https://doi.org/10.18453/rosdok_id00004222

Naranjos Velazquez, N. (2023b). Die Rolle freiberuflicher Hebammen in Netzwerken Frühe Hilfen: Eine quantitative, egozentrierte Netzwerkanalyse. Wiesbaden: Springer VS.

Naranjos Velazquez, N. (2024). ChatGPT als KI-Assistent in der Forschung zu sexuellem Kindesmissbrauch: Wie hoch ist die Übereinstimmung zwischen Mensch und Künstlicher Intelligenz?. Berlin: Jahrestagung DGKIM. https://doi.org/10.13140/RG.2.2.30949.41448

Naranjos Velazquez, N. (in pressa). ChatGPT als KI-Assistent in der Aufbereitung von emotional belastenden Inhalten: Ein Forschungsbericht. In J. Späte, D. C. Stix, Endter, & K. Krauskopf (Eds.), #GesellschaftBilden im Digitalzeitalter (pp. 193-205). Münster: Waxmann.

Naranjos Velazquez, N. (in press[b]). *Adult Disclosure of Childhood Sexual Abuse: Insights into Silent Networks.*

OpenAI. (2023). *GPT-4 Technical Report*. https://arxiv.org/pdf/2303.08774.pdf

OpenAI. (2024). *Best practices for prompt engineering with the OpenAI API.* https://help.openai.com/en/articles/6654000-best-practices-for-prompt-engineering-with-the-openai-api

Perry, B. L., Pescosolido, B. A. & Borgatti, S. P. (2018). *Egocentric network analysis: foundations, methods, and models*. Cambridge University Press.

UKASK. (2024). *Geschichten, die zählen*. https://www.geschichten-die-zaehlen.de

van Manen, M. (2023). What Does ChatGPT Mean for Qualitative Health Research? *Qualitative health research*, *33*(13), pp. 1135–1139. https://doi.org/10.1177/10497323231210816

Williamson, E., Gregory, A., Abrahams, H., Aghtaie, N., Walker, S.-J. & Hester, M. (2020). Secondary Trauma: Emotional Safety in Sensitive Research. *Journal of academic ethics*, *18*(1), pp. 55–70. https://doi.org/10.1007/s10805-019-09348-y

**Contact email:** noreen.naranjos-velazquez@iu.org