## RESEARCH ARTICLE

## USAGE OF LARGE LANGUAGE MODELS FOR ENHANCING INTERACTIVE SYSTEMS: CHALLENGES AND PROSPECTS

**R.A. Cherniavskyi and Y.M. Krainyk**
Petro Mohyla Black Sea National University.

………………………………………………………………………………………………....

| *Manuscript Info* | *Abstract* |
|---|---|
| …………………….. | ……………………………………………………………… |

This study explores the advantages of using Large language models in interactive systems, analyzing existing technical and ethical challenges, and examining development prospects, particularly the optimization of Large language models for local processing on mobile devices to ensure user data privacy and security. Large language models such as GPT-4 and BERT can analyze large volumes of data and understand the context of queries, enabling intuitive interfaces. However, using Large language models involves significant technical and ethical challenges. Technical challenges include optimizing computational resources and ensuring stable model performance on resource-limited mobile devices. Ethical aspects include user data privacy and security issues, as data often needs to be processed on remote servers. This paper proposes solutions to these challenges by optimizing Large language models for local processing on mobile devices.

…………………………………………………………………………………………………....

## Introduction:-

In today's world, interactive technologies are becoming increasingly important, changing the way people interact with digital systems. From voice assistants to smart Internet of Things (IoT) devices, these technologies enable quick and efficient interaction using natural language. The use of large language models (LLMs), such as GPT-4, BERT, and others, is becoming a key factor in ensuring this effective communication. LLMs demonstrate high accuracy in recognizing and processing natural language, allowing for the creation of intuitive interfaces that significantly facilitate user interaction with technology.

Specifically, large language models have the ability to understand the context and intentions of the user, making them indispensable in various applications: from personal assistants to automated customer query processing systems in business. Due to their computational capabilities and ability to learn from vast amounts of data, these models can provide a high level of accuracy and relevance in responses, which is critically important for building effective interactive systems.

The aim of this research is to develop highly efficient solutions to improve interaction between users and technological systems. This involves creating intuitive interfaces based on large language models capable of adequately processing natural language and providing a high level of accuracy and relevance in responses to user

**Corresponding Author:-R.A. Cherniavskyi**
Address**:-**Black Sea National University Named After Peter Mogil.

queries. Additionally, the research focuses on addressing the technical and ethical challenges associated with the use of large language models to ensure data security and privacy.

The development of interactive technologies based on LLMs opens up new opportunities for improving people's quality of life, but at the same time, it poses new challenges for developers and researchers in this field. Solving these challenges requires an interdisciplinary approach that combines knowledge from computer science, ethics, law, and other disciplines. Only in this way can the safe and effective implementation of large language models in our daily lives be ensured.

Analysisofpreviousstudiesshowsthatlargelanguagemodelscansignificantlyenhancetheefficiencyofinteractivesystems. Largelanguagemodelsareideallysuitedfortextclassification,               providinghighaccuracyininterpretingvarioustexts. Inthecontextoftextclassification,          theyhelpdeterminetopics,          categories,          andstyles, simplifyingtheanalysisoflargevolumesofinformation.          Largelanguagemodelsprovide          a reliabletoolforautomatedthematicorientationandidentifyingimportantcharacteristicsoftexts [1].

Studiesalsoemphasizethatlargelanguagemodels,                         suchas                         GPT-4, canunderstandandgeneratehumantextwithhighaccuracy,         buttheiruseisassociatedwithrisksofunwantedbehavior. Forexample,                         modelscanusevariousstrategiestoachievehighscoresfromhumansupervisors, whichcanleadtoundesirableconsequencessuchasdatamanipulationorcircumventingrestrictions [2].

Thestudy             "Thealignmentproblemfrom             a             deeplearningperspective"             [3] analyzesthedifficultiesthatarisewhenattemptingtoalignthebehavioroflargelanguagemodelswithhumanexpectationsand values.                         Theauthorsemphasizethatwithoutsignificanteffortsfromdevelopers, artificialintelligencesystemsmayachievegoalsthatconflictwithhumaninterests.         Thisphenomenon,         knownas "rewardhacking,"         demonstratesthecomplexityofensuringcorrectmodelbehavior.         Oneofthecriticalissuesis "situationalawareness," whichisthemodel'sabilitytounderstandthecontextinwhichitoperatesandadaptitsbehavioraccordingly. Thisincludesusingknowledgeabouttheweaknessesoftrainingalgorithmsandhumansupervisorstoachievesetgoals, evenifthesegoalscontradicttheprimaryintentionsofthedevelopers.                         Overall, previousstudiesshowthattheuseoflargelanguagemodelsininteractivesystemshassignificantpotentialbutrequires         a carefulapproachtotheirdevelopmentandimplementation, consideringtechnicalandethicalaspects.

### Advantages of Using Large Language Models in Interactive Systems:-
Large language models (LLMs) are a type of artificial intelligence based on deep learning aimed at understanding and generating natural language. These models operate on neural networks capable of analyzing and interacting with large volumes of textual information, providing high accuracy and relevance in responses to user queries. The main characteristics of LLMs are deep learning, processing large volumes of data, and contextual understanding.

One of the main advantages of using LLMs in interactive systems is their ability to improve the accuracy and relevance of responses to user queries. Thanks to complex neural networks, such models can analyze the context of queries and provide more precise answers, significantly enhancing user interaction with the system. For example, LLMs can recognize ambiguous words and phrases by understanding the context of their usage, which helps avoid misunderstandings and inaccuracies.

Large language models enable the expansion of the functionality of interactive systems by adding new capabilities for processing queries. For instance, LLMs can be used to automate responses to complex queries that involve multiple stages of data processing. This makes it possible to create more sophisticated and powerful interactive systems that can efficiently handle large volumes of information and ensure high-quality user service.

Specifically, large language models are used to create chatbots that can respond to complex user queries, providing detailed information and helping to solve various problems. Such systems can automatically classify queries, prioritize tasks, and provide relevant responses in real-time.

### TechnicalandEthicalChallenges:
One of the main technical problems faced by large language models (LLMs) is the optimization of computational resources. Large language models require significant computational power for training and deployment, which can

be expensive and challenging. This is especially relevant for mobile devices, where computational resources and energy capabilities are significantly limited. Ensuring the stable operation of such models on mobile devices requires careful optimization of algorithms and efficient use of resources.

Ethical aspects include issues of data privacy and security. Large language models can process vast amounts of personal data, raising concerns about information privacy. In many cases, user data is transferred to remote servers of large companies for processing, which creates the risk of unauthorized use or theft. One possible solution to this problem is optimizing LLMs to work on mobile devices for local processing of user queries. This approach allows data to be stored and processed directly on the user's device, significantly reducing the risks associated with transferring data to remote servers.

Optimizing large language models for mobile devices is a challenging task, but it has significant advantages. Local data processing ensures high privacy and security of user information, as the data does not leave the device. This avoids potential data leaks and reduces dependence on external servers.

Several solutions already exist that allow large language models to run locally on mobile devices. These solutions can be used as a basis for solving the problem of local LLM deployment:

**TensorFlowLite**:
This is a lightweight version of TensorFlow optimized for mobile and embedded devices. TensorFlow Lite supports the execution of machine learning models locally on the device, providing high performance and low latency.

**ONNXRuntimeMobile**:
This is a tool for executing models created in the ONNX (Open Neural Network Exchange) format on mobile devices. ONNX Runtime Mobile ensures high efficiency and allows leveraging the advantages of different hardware platforms.

**CoreML**:
A framework by Apple for deploying machine learning models on iOS devices. Core ML is optimized for working with the limited resources of mobile devices, providing high speed and energy efficiency.

**PyTorchMobile**:
A version of PyTorch optimized for mobile devices. PyTorch Mobile allows deploying and running machine learning models directly on mobile devices, supporting high performance and low latency.

Based on these solutions, an algorithm for optimizing large language models has been developed, allowing for efficient local processing of queries on mobile devices. This ensures high privacy and security of user data, reduces query processing delays, and decreases dependence on internet connections and external servers.

The work focuses on creating local models that can operate on devices with limited resources, using modern machine learning algorithms to reduce computational costs and ensure stable performance. This enhances the efficiency and reliability of interactive systems, providing safe and confidential processing of user data.

**Development Prospects:-**
The development of large language models opens up new opportunities for improving interactive technologies. One area of research is the enhancement of data processing algorithms, which will allow models to provide even more accurate and relevant responses. The use of new machine learning and artificial intelligence methods will improve the accuracy and speed of query processing, making interactive systems even more efficient and reliable. Additionally, the development of new methods to ensure data privacy and security will help mitigate the risks associated with using large language models.

Another important direction is the integration of LLMs with other technologies, such as computer vision and robotics. This integration will allow the creation of complex systems capable of effectively interacting with users and the environment, opening new opportunities for innovation in various fields, from medicine to industry.

The development of data protection and privacy technologies is also critically important. The growing interest in ethical and security issues in the use of LLMs stimulates the development of new solutions for protecting personal data and preventing its unauthorized use. This includes the implementation of more advanced encryption methods, as well as the development of policies and standards that regulate data usage.

Continuous improvement of hardware will also contribute to the development of LLMs. With the emergence of new, more powerful, and energy-efficient processors, it will become possible to deploy even more complex models on mobile devices, expanding their applications and improving efficiency.

Overall, the prospects for the development of large language models in interactive systems include the improvement of algorithms, integration with other technologies, ensuring data privacy, and the advancement of hardware. This will enable the creation of even more efficient, secure, and innovative solutions that enhance user interaction with digital systems.

## Conclusions:-
The use of large language models in interactive systems has significant potential to improve accuracy, efficiency, and overall user experience. However, for their successful implementation, several technical and ethical issues must be addressed. The main technical challenges include optimizing computational resources and ensuring the stable operation of models on mobile devices. Ethical aspects involve issues of user data privacy and security.

Optimizing LLMs for mobile devices allows for local processing of queries, ensuring high privacy and security of user data, reducing query processing delays, and decreasing dependence on external servers. Existing solutions, such as TensorFlow Lite, ONNX Runtime Mobile, Core ML, and PyTorch Mobile, can serve as a basis for creating local models.

Future development prospects include improving data processing algorithms, integrating with other technologies, developing data protection methods, and continuously enhancing hardware. This will enable the creation of more efficient, secure, and innovative solutions that enhance user interaction with digital systems.

## References:-
1. Kann E., Cotterell R., Gonen M., Gorman J., Zhao Y., Liu S., Palmer M., Goyal A. NorthAmericanChapteroftheAssociationforComputationalLinguistics (NAACL) 2024: Tutorials // Proceedingsofthe 2024 ConferenceoftheNorthAmericanChapteroftheAssociationforComputationalLinguistics: Tutorials. NewYork, NY, USA, 2024. C. 1-27.
2. Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., ... Le, Q. V. (2022). EmergentAbilitiesofLargeLanguageModels.
3. Ngo R., Chan L., Mindermann S. Thealignmentproblemfrom a deeplearningperspective.