# Summer School 2024 "Natural Language Processing" 27 May to 8 June 2024

The Summer School 2024 on "Natural Language Processing" organised by International Centre for Free and Open Source Solutions (ICFOSS) was held from 27 May to 08 June 2024 at ICFOSS. Targeting post-graduate students, research scholars, faculty members, and industry professionals, the Summer School featured project-oriented sessions that combined practical concepts with theoretical knowledge. This approach provided participants with a comprehensive understanding of Natural Language Processing (NLP), fostering a collaborative and inclusive learning environment.



# **Inaugural Ceremony**

The Summer School 2024 on "Natural Language Processing" commenced on 27 May 2024 with an inaugural ceremony. The event began with a warm welcome speech by Dr. Rajeev R R , Program Head and Associate Professor at ICFOSS. The formal inauguration was carried out by Dr. Sunil TT, the esteemed Director of ICFOSS. Following the inauguration, Ms. Chithra M S, Secretary and Registrar of ICFOSS delivered the presidential address. The event then moved forward with the keynote address by Dr. Sunil TT, where he elaborated on the current trends and future prospects of NLP. The Summer School promises to be highly beneficial for the participants, offering a rich program that includes lectures, hands-on sessions, and interactive discussions. Participants will have the opportunity to learn from leading experts in the field, engage in collaborative projects, and gain practical experience with the NLP tools and techniques.





# Day 1

### Resource Person : Dr. Rajeev RR

Topic : Introduction to Natural Language Processing (NLP)

Time : 2 PM to 5 PM

Following the inaugural ceremony, Dr Rajeev RR took a session on "Introduction to NLP". He provided a comprehensive overview of the field, highlighting its significance and applications in today's world. The session aimed to familiarise participants with the basic concepts and techniques used in NLP. He began by defining NLP as a branch of artificial intelligence that focuses on enabling computers to understand, interpret, and generate human language. He explained the importance of NLP in various real-world applications, such as chatbots, language translation, sentiment analysis, and information extraction. The session covered several basic techniques used in NLP, including tokenization, which involves breaking text into individual words or tokens; part-of-

speech tagging, which assigns grammatical tags to words; and syntactic parsing, which analyses the grammatical structure of sentences. He also discussed the wide range of applications of NLP in today's world, including machine translation, speech recognition, text summarization, and sentiment analysis. He emphasised how NLP has revolutionised the way we interact with technology, enabling more natural and intuitive communication. The session concluded with a discussion on the challenges faced by NLP, such as ambiguity in language, understanding context, and handling different languages and dialects. Dr. Rajeev RR also highlighted the future directions of NLP, including the integration of NLP with other technologies such as machine learning and deep learning to further enhance its capabilities.



The session provided a comprehensive overview of the field, covering its basic concepts, techniques, applications, and future directions. The session was insightful and informative, providing participants with a solid foundation in NLP and its potential impact on society.



# Day 2

### Resource Person : Dr. Umesh P

**Topic :** Mathematical Foundation for Deep Learning

Time: 10 AM to 5 PM

Dr. Umesh P's session provided a comprehensive view on the mathematical foundation for deep learning, offering attendees a detailed overview of key concepts essential for understanding and applying deep learning algorithms in practice.

The session began with an overview of machine learning, highlighting its various types, including supervised, unsupervised, and reinforcement learning. Supervised learning involves learning from labelled training data, where each example is paired with the correct output, aiming to learn a mapping from inputs to outputs. Unsupervised learning deals with learning from unlabeled data, seeking to find patterns or structure in the data without explicit guidance. Reinforcement learning centres on learning how to make sequences of decisions based on feedback in the form of rewards or punishments.



Linear algebra's role in machine learning was explained, emphasising its significance in understanding and manipulating high-dimensional data. Vectors are used to represent data points, while matrices are used to represent datasets or transformations. Operations such as matrix multiplication and transpose were underscored for their importance in various machine learning algorithms.

Calculus was discussed as essential for optimization, which involves finding the best parameters for a model to minimise a loss function. Derivatives and gradients were explained as crucial tools for updating parameters in machine learning models, with the Jacobian matrix and Hessian matrix being instrumental in computing higher-order derivatives.

Probability and statistics were highlighted as crucial for modelling

uncertainty in data and making probabilistic predictions. Concepts such as mean, median, mode, covariance, variance, and correlation coefficient were explained as fundamental measures for describing and analysing data distributions. Probability distributions, conditional probability, and probability mass function were explained as essential for modelling uncertainty.

The session also gave an introduction to neural networks, the building blocks of deep learning. Linear regression and logistic regression were discussed as foundational models, with activation functions introducing non-linearity to neural networks. Gradient descent was detailed as an optimization algorithm used to minimise the loss function of a neural network, with regularisation techniques being employed to prevent overfitting. Classification was discussed as a fundamental task in machine learning, with learning algorithm selection being emphasised as crucial for the success of a machine learning project. Model performance assessment metrics such as accuracy, precision, recall, and F1 score were also discussed. The session also covered the feedforward and backpropagation operations in neural networks, providing a detailed explanation of these fundamental processes. Feedforward operation involves the propagation of input data through the neural network, passing through multiple layers of neurons and applying activation functions to produce an output. Each neuron in the network computes a weighted sum of its inputs, adds a bias term, and applies an activation function to produce an output.



Backpropagation operation, on the other hand, is used to update the weights and biases of the neural network based on the error between the predicted output and the actual output. It involves calculating the gradient of the loss function with respect to the weights and biases of the network, and then using this gradient to update the weights and biases using an optimization algorithm such as gradient descent.

During the session, a worksheet on neural networks was provided, which included concepts such as bias, learning rate, cost function, loss function, and the difference between loss and cost functions.

In conclusion, Dr. Umesh P's session provided a comprehensive understanding of the mathematical foundation for deep learning, equipping attendees with the knowledge necessary to apply deep learning algorithms effectively in real-world scenarios.

## Day 3

### Resource Person : Dr. Selvaraj R and Mr. Jibin Kiran

Topic : Linguistic Aspects of NLP

Time: 10 AM to 5 PM

The session on the Linguistic Aspects of Natural Language Processing (NLP) focused on the intricate relationship between language and technology. They began by defining NLP as a branch of artificial intelligence that focuses on enabling computers to understand, interpret, and generate human language. Central to NLP is the concept of "text," which refers to any form of written or spoken language that computers analyse. They emphasised the importance of linguistic cohesion, which involves the grammatical and lexical relationships that contribute to the overall coherence of a text.

The session then explored various tasks involved in NLP. For instance, a "sentence splitter" identifies and separates sentences within a text, while a "tokenizer" breaks down the text into individual tokens, such as words or punctuation marks. Additionally, a "morphological analyzer" examines the structure of words to determine their grammatical properties, and a "POS tagger" assigns grammatical tags to words based on their parts of speech. These tasks are fundamental in enabling computers to analyse and understand the structure of human language.



They also discussed the processing of a text, which involves several stages such as preprocessing, deep analysis, and shallow analysis. Preprocessing includes tasks like cleaning and formatting raw text for analysis, while deep analysis involves in-depth linguistic analysis, including syntactic and semantic analysis. Shallow analysis, on the other hand, focuses on surfacelevel linguistic features.

Furthermore, the session covered computational requirements for NLP tasks, highlighting the need for significant computational resources, especially for tasks requiring deep linguistic analysis. The complexity of NLP tasks and the size of the text being analysed can impact the computational requirements.



In conclusion, the session provided a detailed overview of the linguistic aspects of NLP, showcasing how computational techniques and linguistic knowledge are integrated to enable computers to interact with human language effectively.

# Day 4

Resource Person : Dr. Gopakumar G

Topic : Introduction to Neural Networks, RNN, LSTM

Time: 10 AM to 5 PM

Dr. Gopakumar G's session on the introduction to Neural Networks, Recurrent Neural Networks (RNNs), and Long Short-Term Memory (LSTM) networks provided a comprehensive overview of these fundamental concepts in deep learning. The session began with an introduction to Neural Networks, highlighting their ability to learn complex patterns in data through interconnected layers of neurons. The concept of Multilayer Neural Networks was discussed, demonstrating how multiple layers allow for more sophisticated modelling of data.



The session then focused on the feedforward operation and classification in Neural Networks, explaining how input data is processed through the network to produce an output, which is then used for classification or prediction tasks. The Backpropagation Algorithm was detailed as a method for updating the weights of the network based on the error between the predicted and actual output, allowing the network to learn from its mistakes and improve its performance.

Two modes of operation were discussed: feedforward and learning. In feedforward mode, the network processes input data to produce an output, while in learning mode, the network adjusts its weights based on the error in the output to improve its performance. The session also covered network learning, including learning curves, which depict the network's performance over time, and error surfaces, which represent the relationship between the network's weights and its error.

Backpropagation was explained as a form of feature mapping, where the network learns to map input features to output predictions. Practical techniques for improving backpropagation were discussed, including the choice of activation function, scaling input and target values, training with noise, and manufacturing data to augment the training set.

Other factors influencing backpropagation performance, such as the number of hidden units, initialising weights, learning rates, momentum, weight decay, and criteria function, were also covered. The session highlighted the importance of choosing the right training method (on-line, stochastic, or batch training) and knowing when to stop training to avoid overfitting. The number of hidden layers was discussed as a factor influencing the network's ability to learn complex patterns.

In addition to Neural Networks, the session covered Recurrent Neural Networks (RNNs), which are designed to handle sequential data by maintaining a hidden state that captures information from previous time steps. The structure of RNNs was explained, illustrating how the hidden state is updated at each time step using an activation function. The session also discussed the different categories of sequence modelling, including one-to-one, one-to-many, many-to-one, and many-to-many models, showcasing the versatility of RNNs in various tasks.

The power of RNNs in representation was emphasised, highlighting their ability to capture long-range dependencies in sequential data. Training RNNs was explained, focusing on the Backpropagation through Time (BPTT) algorithm, which extends backpropagation to sequential data by unfolding the network over time. The problem of Vanishing gradients in RNNs was addressed, and Truncated BPTT was introduced as a technique to mitigate this issue by limiting the number of time steps considered during training.



Long Short-Term Memory (LSTM) networks were introduced as a variant of RNNs designed to better capture long-range dependencies. The forward pass of LSTM was explained, detailing how it handles input, forget, and output gates to control the flow of information through the network. Training LSTM was discussed, highlighting the importance of initialising weights and choosing appropriate hyperparameters.

Gated Recurrent Unit (GRU) was introduced as another variant of RNNs, which simplifies the architecture of LSTM while achieving comparable performance. Bidirectional RNNs were discussed as a way to improve the representation of sequential data by processing it in both forward and backward directions. Practical applications of RNNs were explored, including machine translation, where RNNs are used to translate text from one language to another, and generating image descriptions, where RNNs generate textual descriptions of images. The computational graph of RNNs was illustrated, showcasing how data flows through the network during training and inference. Sequence-to-sequence models were discussed as a general framework for tasks that involve mapping input sequences to output sequences, such as machine translation.

Overall, Dr. Gopakumar G's session provided a comprehensive understanding of Neural Networks, RNNs, and LSTMs, equipping attendees with the knowledge necessary to apply these concepts to a wide range of tasks in deep learning.

## Day 5

Resource Person : Ms. Meharuniza Nazeem

**Topic :** Text Processing Basics, Text Representation, Word Embeddings

Time: 10 AM to 1 PM

Meharuniza Nazeem's session on word embedding covered three popular techniques: Word2Vec, FastText, GloVe and One-Hot Encoding. Word embedding is a technique used in natural language processing (NLP) to represent words as vectors in a continuous vector space. These word vectors capture semantic relationships between words, enabling NLP models to better understand and process language.



Word2Vec is a neural network-based model that learns word embeddings by predicting the surrounding words in a sentence (skip-gram model) or predicting a word based on its context (continuous bag of words model). The resulting word vectors capture semantic relationships, such as similarity and analogy, between words. Word2Vec has been widely used in various NLP tasks, including sentiment analysis, machine translation, and named entity recognition, due to its effectiveness in capturing word semantics.

FastText is an extension of Word2Vec that introduces subword information into word embeddings. Instead of representing each word as a single vector, FastText represents words as the sum of their character n-grams. This allows FastText to capture morphological information and handle out-of-vocabulary words more effectively than Word2Vec. FastText has been particularly useful for tasks involving morphologically rich languages or domains with limited training data.

GloVe is another popular word embedding technique that learns word vectors by factoring a matrix of word co-occurrence statistics. UnlikeWord2Vec and FastText, which are based on neural networks, GloVe is based on matrix factorization and directly optimises word vectors to capture global word-word co-occurrence statistics. This results in word embeddings that are effective at capturing global semantic relationships between words.

One-hot encoding is a technique to represent categorical data, such as words, as binary vectors. Each word is represented by a vector where all elements are 0 except for one element, which is 1 at the index corresponding to the word's position in the vocabulary. While simple, one-hot encoding does not capture semantic relationships between words and results in high-dimensional sparse vectors.

Word embedding techniques like Word2Vec, FastText, and GloVe are widely used in NLP for tasks such as sentiment analysis, machine translation, and named entity recognition. These techniques enable machine learning models to better understand and process natural language by representing words in a continuous vector space.

Meharuniza Nazeem's session provided a comprehensive overview of Word2Vec, FastText, and GloVe, highlighting their strengths and applications in NLP. These word embedding techniques have significantly advanced the field of NLP by enabling models to effectively capture semantic relationships between words and improve performance on various NLP tasks. Understanding these techniques is crucial for researchers and practitioners in the field of NLP to develop more accurate and efficient language processing systems.

#### Afternoon session

#### **Resource person :** Mr. Arun A

Topic : Machine Learning through NLP

#### Time: 2 PM to 5 PM

Mr. Arun A's session on Machine Learning through Natural Language Processing (NLP) provided a detailed exploration of how machine learning techniques are applied in various aspects of NLP. He began by discussing the different types of machine learning, including supervised, unsupervised, and deep learning, and their applications in NLP tasks such as text classification, sentiment analysis, and named entity recognition. He highlighted the importance of feature engineering in machine learning for NLP, where linguistic features such as word embeddings and syntactic structures are used to train models effectively.

The session also covered the challenges of working with large-scale datasets in NLP and the need for scalable machine learning algorithms. He discussed the role of neural networks in NLP, including convolutional neural networks (CNNs) for text classification and recurrent neural networks (RNNs) for sequence modelling. He also highlighted the significance of attention mechanisms in improving the performance of machine learning models in NLP tasks such as machine translation and text summarization.



Moreover, He addressed the ethical considerations in using machine learning in NLP, emphasising the need for fairness, transparency, and accountability in AI systems. He discussed the importance of bias detection and mitigation strategies to ensure that machine learning models are not perpetuating or amplifying biases present in the data.

In conclusion, Mr. Arun A session provided a comprehensive overview of machine learning through NLP, showcasing its potential to revolutionise language understanding and interaction. His insights into the latest advancements and challenges in the field highlighted the need for continued research and innovation to harness the full power of machine learning in NLP.

# Day 6

### **Resource Person : Mr. Sabeerali KP**

Topic : Introduction to Transformers - Understanding the transformer

architecture, Overview of Attention Mechanism

### Time: 10 AM to 1 PM

The objective of this session was to provide participants with an in-depth understanding of Transformers, a revolutionary deep learning architecture that has transformed the field of Natural Language Processing (NLP). The session aimed to elucidate the architecture, working principles, and applications of Transformers, along with recent advancements and future directions.



The session began with an introduction to Transformers, a deep learning architecture introduced in the paper "Attention is All You Need" by Vaswani et al. in 2017. Transformers have gained widespread popularity in NLP tasks due to their ability to model long-range dependencies efficiently and their parallelizable nature, enabling faster training and inference compared to recurrent neural networks (RNNs) and convolutional neural networks (CNNs).

Participants were introduced to the key components of Transformers, including:

- 1. Self-Attention Mechanism: Transformers rely on self-attention mechanisms to weigh the importance of different words in a sequence when generating representations. This mechanism enables the model to capture global dependencies and contextual information effectively.
- Multi-Head Attention: To enhance representational capacity, Transformers employ multi-head attention mechanisms, allowing the model to focus on different parts of the input sequence simultaneously.
- 3. Positional Encoding: Since Transformers lack recurrence and convolution, they require positional information to understand the order of words in a sequence. Positional encodings are added to the input embeddings to provide this information.
- 4. Feedforward Neural Networks: Transformers include feedforward neural networks to process the representations generated by the self-attention mechanism and produce the final output.

The session covered various applications of Transformers across NLP tasks, including:

- 1. Language Translation: Transformers, particularly the Transformer model introduced in the original paper, have been highly successful in machine translation tasks, achieving state-of-the-art performance on benchmark datasets such as WMT.
- 2. Question Answering: Transformers have been applied to question

answering tasks, where they excel at processing and understanding contextual information to generate accurate answers.

- 3. Text Summarization: Transformers are effective in text summarization tasks, where they can generate concise summaries of longer documents by focusing on important information.
- 4. Named Entity Recognition (NER): Transformers have been applied to NER tasks, where they can identify and classify named entities such as names of people, organisations, or locations.

Participants were briefed on recent advancements in Transformer architecture, including:

- 1. BERT (Bidirectional Encoder Representations from Transformers): BERT, introduced by Devlin et al. in 2018, is a pre-trained Transformer model that has achieved significant improvements in various NLP tasks by pre-training on large corpora of text data.
- 2. GPT (Generative Pre-trained Transformer): GPT models, developed by OpenAI, are autoregressive Transformer models capable of generating coherent and contextually relevant text.
- 3. XLNet: XLNet, introduced by Yang et al. in 2019, is a generalised autoregressive pre-training method that outperforms previous Transformer-based models by leveraging permutation-based training objectives.

In conclusion, the session provided participants with a comprehensive understanding of Transformers, including their architecture, working principles, applications, recent advancements, and future directions. By grasping the fundamentals of Transformers, participants are better equipped to leverage this powerful deep learning architecture in various NLP tasks and contribute to advancements in the field.

#### **Afternoon Session**

Resource Person : Mr. Sabeerali KP

Topic : Deep Learning Frameworks

Time: 10 AM to 5 PM

The session on Day 6, conducted by Mr. Sabeerali KP, focused on Deep Learning Frameworks, with a detailed discussion on major frameworks like TensorFlow and PyTorch. Participants were engaged in a hands-on session, where they actively built models using these cutting-edge tools.

The morning session specifically covered TensorFlow, an open-source deep learning framework developed by Google in 2015. TensorFlow is renowned for its extensive documentation, training support, scalable production and deployment options, multiple abstraction levels, and compatibility with various platforms, including Android. It is a symbolic maths library ideal for neural networks and excels in dataflow programming for diverse tasks.

During the hands-on segment, participants were guided through the process of creating a model using the FashionMNIST dataset. The session started with importing the necessary libraries and datasets, followed by data preprocessing, model building, training, and evaluation. Participants learned how to define a neural network architecture using TensorFlow's high-level Keras API, compile the model with an optimizer and loss function, and train it using the training data. They also explored techniques

for model evaluation and learned how to make predictions on new data. Overall, the hands-on session provided participants with practical experience in using TensorFlow to build and train deep learning models.

Then focused on PyTorch, a machine learning framework based on the Torch library, originally developed by Meta AI and now part of the Linux Foundation umbrella. PyTorch is widely used for applications such as computer vision and natural language processing. It is known for its dynamic computational graph, which allows for more flexible and intuitive model building compared to static graph frameworks like TensorFlow. PyTorch is free and open-source software released under the modified BSD licence.

During the session, participants built the same model using the FashionMNIST dataset in PyTorch to gain a deeper understanding of the differences between PyTorch and TensorFlow. They imported the required libraries and datasets using specific commands, similar to the process in TensorFlow. By comparing the two frameworks in the context of the same task, participants were able to appreciate the unique features and advantages of each framework, enabling them to make informed decisions in choosing the right framework for their future projects.

## Day 7

Resource Person : Dr. Nisha Varghese

**Topic :** Advanced NLP Techniques-Transfer learning in NLP, Introduction to pre-trained language models (e.g., BERT)

### Time: 10 AM to 5 PM

The session commenced with an introduction to the concept of transfer learning in Natural Language Processing (NLP), where Dr. Nisha elaborated on how pre-trained models can be adapted to new, related tasks. She emphasised the efficiency of transfer learning in terms of reducing data and computational requirements while enhancing model performance.



Dr. Nisha then elaborated the details of pre-trained language models, explaining their role in revolutionising NLP. She highlighted the evolution of these models, starting from simpler word embeddings like Word2Vec and GloVe to more sophisticated models like ELMo, GPT, and BERT. The focus was primarily on BERT (Bidirectional Encoder Representations from Transformers), a model known for its bidirectional training approach that allows it to understand the context of a word based on its surrounding

words. Dr. Nisha detailed BERT's architecture, explaining how it uses transformers to process text in a way that captures the intricate relationships between words in a sentence.



One of the key aspects discussed was the process of pre-training and finetuning BERT. Dr. Nisha explained that BERT is pre-trained on a massive corpus of text data, allowing it to learn a wide range of linguistic features. This pre-trained model can then be fine-tuned on specific downstream tasks, such as text classification, question answering, and sentiment analysis, with much smaller datasets. She demonstrated how this approach significantly reduces the time and computational resources needed compared to training a model from scratch.

Dr. Nisha also highlighted the numerous advantages of using pre-trained models. These models save time and resources while achieving state-ofthe-art results on many NLP benchmarks. They are also highly versatile, capable of being adapted to various NLP tasks with minimal adjustments. This versatility makes them invaluable tools for both research and practical applications in NLP.



To provide a hands-on experience, Dr. Nisha conducted a practical demonstration on fine-tuning BERT for a specific NLP task using the Hugging Face Transformers library. Participants were guided through the process of setting up the environment, preparing data, and running the model. This step-by-step guidance made the complex process accessible and manageable, allowing participants to gain practical experience and confidence in implementing these techniques.

The session concluded with an interactive Q&A segment, where Dr. Nisha addressed various queries from the participants. She provided clarity on the implementation challenges and shared insights on the future directions of NLP and the potential advancements in pre-trained models. The interactive Q&A session further enriched the learning experience, with discussions on practical challenges and future directions in NLP.

Overall, the session by Dr. Nisha Varghese was highly informative and well-received by the participants. It provided a comprehensive understanding of advanced NLP techniques, particularly the significance of transfer learning and the practical application of pre-trained language models like BERT. The hands-on demonstration was particularly beneficial, allowing participants to gain practical experience in implementing these techniques. The session left participants better equipped to utilise these technologies in their work, enhancing their capabilities in the rapidly evolving field of NLP.

## Day 8

Resource Person : Dr. Nisha Varghese

**Topic :** Auto-regressive models

Time: 10 AM to 5 PM

The session began with an introduction to the concept of auto-regressive models. Dr. Nisha explained that auto-regressive models predict future values based on previously observed values, a technique that has proven particularly effective in language modelling. These models generate text by predicting the next word in a sequence, given the preceding words, making them essential tools for tasks such as text generation, translation, and completion.



Dr. Nisha elaborated on the foundational principles of auto-regressive models, starting with simpler models like n-grams and moving towards more sophisticated models like the Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM) networks, and Transformers. She highlighted the limitations of traditional n-gram models, such as their inability to capture long-range dependencies, which paved the way for the development of RNNs and LSTMs. These models, with their ability to maintain context over longer sequences, significantly improved language modelling capabilities.

The session then focused on the Transformer architecture, which has become the backbone of modern auto-regressive models. Dr. Nisha explained how Transformers, with their self-attention mechanisms, address the shortcomings of RNNs and LSTMs by enabling parallel processing and better handling of long-range dependencies. She provided a detailed overview of popular auto-regressive models built on the Transformer architecture, such as GPT (Generative Pre-trained Transformer) and its variants.

One of the key points discussed was the training process of these models. Dr. Nisha explained the pre-training and fine-tuning steps, similar to those in transfer learning. In pre-training, the model learns general language features from a large corpus of text. During fine-tuning, the model is adapted to specific tasks with smaller, task-specific datasets. This two-step process allows auto-regressive models to achieve high performance on various NLP tasks with relatively little task-specific data.

Dr. Nisha also highlighted the practical applications of auto-regressive models. These models are used in a wide range of NLP tasks, including text generation, machine translation, dialogue systems, and more. She demonstrated how auto-regressive models can be employed to generate coherent and contextually relevant text, showcasing examples from applications like chatbots and automated content creation.

To provide a practical understanding, Dr. Nisha conducted a hands-on demonstration. Participants were guided through the process of implementing an auto-regressive model using the Hugging Face Transformers library. This demonstration included setting up the environment, preparing data, and running the model. The step-by-step guidance helped participants grasp the implementation details and gain practical experience in using these models. The session concluded with an interactive Q&A segment. Dr. Nisha addressed various queries from the participants, providing further insights into the challenges and nuances of implementing auto-regressive models. She also discussed the future directions of NLP, emphasising the continuous evolution and potential





The interactive Q&A session further enriched the learning experience, fostering a deeper understanding of the practical applications and future potential of auto-regressive models in NLP.

Overall, Dr. Nisha Varghese's session on auto-regressive models was highly informative and well-received. It provided participants with a thorough understanding of the principles, applications, and practical implementation of auto-regressive models in NLP. The session left participants better equipped to utilize these models in their work, enhancing their capabilities in the rapidly evolving field of NLP.

## Day 9

**Resource Person :** Ms. Meharuniza Nazeem

Topic : Explainable AI in NLP

### Time: 10 AM to 1 PM

Ms. Meharuniza Nazeem took an informative session on "Explainable AI in Machine Learning" elucidating the concepts and significance of explainable artificial intelligence (XAI) in the context of machine learning. The session commenced with an overview of the growing importance of XAI, particularly in fields where transparency and interpretability of AI models are critical, such as healthcare, finance, and legal sectors.



Ms. Meharuniza began by defining explainable AI and its role in making machine learning models understandable to humans. She highlighted the challenges posed by black-box models, where decisions are opaque and difficult to interpret, contrasting this with the need for models that can provide insights into their decision-making processes. The session explored various techniques and methodologies used in XAI, including feature importance methods, surrogate models, and local explanations. Ms. Meharuniza explained how these techniques aim to uncover the internal workings of AI models, allowing stakeholders to understand how inputs are processed and decisions are made. She emphasised the balance between model complexity and explainability, illustrating how simpler models often provide clearer explanations but may sacrifice predictive accuracy.

Key concepts such as model-agnostic techniques (e.g., SHAP, LIME) and model-specific approaches (e.g., decision trees, rule extraction) were discussed in detail. Ms. Meharuniza provided practical examples and case studies where XAI has been successfully applied to improve model transparency and user trust.

Participants were engaged in interactive discussions and demonstrations, where they explored real-world applications of XAI across different domains. Ms. Meharuniza addressed participant queries, providing insights into the implementation challenges and best practices for integrating XAI into machine learning pipelines.

The session concluded with a forward-looking discussion on the future of XAI, highlighting ongoing research trends and the evolving regulatory landscape around AI transparency. Ms. Meharuniza encouraged participants to consider ethical implications and societal impacts when developing and deploying AI systems.

Overall, Ms. Meharuniza Nazeem's session on "Explainable AI in Machine Learning" was well-received and highly beneficial for participants seeking a deeper understanding of how XAI enhances transparency and trust in AI systems. The session provided practical insights and tools that participants can utilise to enhance the explainability of their own machine learning models, thereby advancing the responsible adoption of AI technologies in various applications.

# **Day 10**

Resource Person : Mr. Navaneeth S and Mr. Arun A

**Topic :** Hugging face, nanoGPT, GPT for all

Time: 10 AM to 5 PM

The session covered the Hugging Face library, a popular open-source library for natural language processing tasks. Hugging Face provides a wide range of pre-trained models, including BERT, GPT, and many others, along with easy-to-use interfaces for model loading, fine-tuning, and inference.

Participants were provided with practical demonstrations of using BERT and the Hugging Face library for various NLP tasks, including:

- 1. Text Classification: Using pre-trained BERT models for text classification tasks, such as sentiment analysis or topic classification.
- 2. Named Entity Recognition (NER): Fine-tuning BERT for NER tasks to identify and classify named entities in text, such as names of people, organisations, or locations.
- 3. Question Answering: Using BERT for question answering tasks, where the model is tasked with providing answers to questions based on given contexts.



The session also covered advanced features of the Hugging Face library, including:

- 1. Pipeline API: Hugging Face provides a pipeline API that allows users to perform various NLP tasks such as text generation, translation, and summarization with pre-trained models in a few lines of code.
- 2. Model Hub: Hugging Face's Model Hub provides a repository of pre-trained models and community-contributed models that can be easily accessed and used for different tasks.
- 3. Tokenizer: Hugging Face provides tokenizers for various pre-trained models, allowing users to tokenize text data consistently across different models and frameworks.

The objective of this session was to provide participants with an in-depth understanding of NanoGPT, a lightweight and efficient variant of the Generative Pre-trained Transformer (GPT) model developed by OpenAI. The session aimed to elucidate the architecture, capabilities, and potential applications of NanoGPT, as well as its significance in democratising access to powerful language models.



The session commenced with an introduction to NanoGPT, a scaled-down version of the GPT model designed for resource-constrained environments such as mobile devices, edge computing devices, and IoT devices. NanoGPT retains the core architecture and capabilities of GPT while significantly reducing the model size and computational requirements, making it suitable for deployment on devices with limited memory and processing power.

Participants were introduced to the key features of NanoGPT, including:

- 1. Lightweight Architecture: NanoGPT utilises a compact architecture with fewer parameters compared to its larger counterparts, enabling efficient inference on resource-constrained devices.
- 2. Efficient Inference: NanoGPT is optimised for inference speed and

memory footprint, allowing it to generate text responses quickly and without excessive resource consumption.

- 3. Scalability: Despite its reduced size, NanoGPT retains scalability and can be fine-tuned on domain-specific data to adapt to specific tasks and applications.
- 4. Language Understanding: NanoGPT can generate coherent and contextually relevant text responses, making it suitable for applications such as chatbots, virtual assistants, and text generation tasks.

The session covered various potential applications of NanoGPT across different domains, including:

- 5. Conversational AI: NanoGPT can be deployed in chatbots, virtual assistants, and conversational agents to provide natural and engaging interactions with users.
- 6. Text Generation: NanoGPT can generate creative and contextually relevant text in applications such as content creation, storytelling, and dialogue generation.
- 7. Language Translation: NanoGPT can be used for real-time language translation on mobile devices, enabling offline translation capabilities without relying on cloud-based services.
- 8. Personalization: NanoGPT can be fine-tuned on user-specific data to provide personalised recommendations, responses, and experiences in applications such as content curation and recommendation systems.



The session also addressed challenges and considerations associated with deploying NanoGPT on resource-constrained devices, including:

- 1. Model Size vs. Performance Trade-off: Balancing model size with performance and accuracy is crucial when deploying NanoGPT on devices with limited resources.
- 2. Optimization Techniques: Various optimization techniques such as quantization, pruning, and knowledge distillation can be employed to further reduce the size and computational requirements of NanoGPT.
- 3. Privacy and Security: Ensuring privacy and security of user data when deploying NanoGPT on edge devices is essential, requiring robust encryption and data protection measures. The session provided participants with a comprehensive understanding of NanoGPT, including its architecture, features, potential applications, and challenges. By grasping the fundamentals of NanoGPT, participants got a better idea to utilise this lightweight and efficient

language model for various applications on resource-constrained devices, contributing to advancements in edge computing and democratising access to AI technologies.

## **Day 11**

### Resource Person : Mr. Arun A

**Topic :** Future Trends and Project Ideas-Emerging trends in NLP & Artificial Intelligence, Brainstorming and discussion on potential NLP projects

Time: 10 AM to 1 PM

Mr. Arun A led an engaging session on "Future Trends and Project Ideas: Emerging Trends in NLP & Artificial Intelligence". The session aimed to explore current and upcoming trends in Natural Language Processing (NLP) and Artificial Intelligence (AI), while facilitating brainstorming and discussions on potential NLP projects.

The session commenced with an overview of recent advancements and emerging trends in NLP and AI technologies. He discussed the evolution from traditional machine learning approaches to deep learning techniques, emphasising the role of large-scale language models, such as BERT and GPT, in revolutionising NLP tasks like language understanding, generation, and translation.

Key topics covered included the integration of multimodal AI, which combines textual and visual data for more robust understanding and interaction. He highlighted the applications in image captioning, video summarization, and multimodal sentiment analysis, illustrating how these advancements are shaping future AI applications.

He included a brainstorming segment where participants actively contributed project ideas leveraging these emerging trends. He facilitated discussions on the feasibility, scope, and impact of various project proposals, encouraging participants to consider real-world challenges and opportunities in implementing NLP solutions.

Participants engaged in group activities to refine project concepts, exploring potential datasets, methodologies, and evaluation metrics. The session emphasised the importance of interdisciplinary collaboration, drawing on expertise from linguistics, computer science, and domainspecific knowledge to tackle complex NLP problems effectively.

He shared case studies and success stories from industry and research, demonstrating how innovative NLP projects have addressed societal needs, improved business processes, and enhanced user experiences. He encouraged participants to think creatively and critically about the ethical implications and societal impacts of their proposed projects.

The session concluded with a summary of key insights and actionable takeaways for participants interested in pursuing NLP projects. He underscored the importance of continuous learning and adaptation to keep pace with rapidly evolving NLP technologies and trends.

Overall, the session on "Future Trends and Project Ideas: Emerging Trends in NLP & Artificial Intelligence" by Mr Arun A was instrumental in fostering a collaborative environment for exploring innovative NLP solutions. It equipped participants with valuable knowledge, inspiration, and practical strategies to embark on impactful NLP projects, contributing to the advancement of AI technologies and applications in diverse domains.

# **Day 12**

Day 12 of the summer school commenced with the project presentations by the participants. The project "Text Summarization," guided by Ms. Parvathy Raj, secured first place. The project "Text Generation" guided by Mr. Muhammed Fawzan secured second place. The afternoon was dedicated to the valedictory ceremony, where students received certificates and gifts for their achievements and participation throughout the program.

# **Cultural Events**

Cultural events were also conducted in between the summer school. These events provided a refreshing break from the intensive academic sessions and allowed participants to unwind, socialise, and experience the warmth of music and dance. Participants performed various programs during the event, showcasing their talents and creating a vibrant, energetic atmosphere. The events were meticulously organized, with dedicated time slots in the evenings, ensuring they did not interfere with the academic schedule These activities fostered a sense of community and camaraderie, enhancing the overall experience of the summer school. The events were not only a platform for relaxation but also an opportunity for cultural exchange, allowing participants to appreciate and celebrate the diversity among them. Overall, the cultural events were a vital aspect of the summer school, providing a balanced experience that combined rigorous academic learning with enriching cultural interactions.



# **Valedictory Ceremony**

The 12-day Summer School on "Natural Language Processing" concluded with a Valedictory Ceremony. The valedictory ceremony was presided over by Dr. Rajeev RR (Head e-Governance & Development, ICFOSS), who delivered the closing remarks, appreciating the efforts of the faculty members and all contributors to the program's success. Participants then shared their valuable feedback, which will be instrumental in enhancing and improving the summer school for future participants. Dr. Rajeev RR distributed certificates to participants for their active engagement throughout the program. Awards were given to the best projects, with the "Text Summarization" project, mentored by Parvathy Raj, securing first place, and the "Text Generation" project, mentored by Mr. Muhammed Fawzan A, securing second place. Additionally, speakers, project guides, and assistants who contributed to the program's success were duly

### acknowledged.





# **Event Outcome**

The ICFOSS Summer School 2024 on Natural Language Processing (NLP) concluded with great success, marked by the keen enthusiasm and active participation of attendees. The participants, who were divided into

groups, were particularly keen on presenting their project proposals. This competitive spirit culminated in the best project from each group winning prizes, recognizing their exceptional contributions and innovative ideas in the field of NLP.

The event focused heavily on the academic and research development of Natural Language Processing applications. Throughout the program, participants engaged in intensive learning sessions, workshops, and handson projects that allowed them to explore and apply advanced NLP techniques. The sessions, led by esteemed experts in the field, covered a range of topics from the basics to the latest advancements in NLP.

In addition to the academic rigour, the program also emphasised practical application. Participants worked on real-world problems, applying their newly acquired knowledge to develop innovative solutions. These projects not only showcased their technical skills but also their ability to think critically and creatively.

The cultural events conducted during the summer school provided a refreshing break from the intensive academic schedule. Participants showcased their talents in various programs fostering a sense of community and camaraderie. The program concluded with the project presentations by the participants. A panel of judges, consisting of experts, evaluated the projects based on innovation, practical application, and presentation. The winning projects received well-deserved recognition and prizes, adding a competitive edge and motivating participants to strive for excellence.

Overall, the Summer School 2024 on Natural Language Processing was a

resounding success. It provided participants with valuable knowledge, practical experience, and a platform to showcase their talents and innovations.