1.  **What are Vanilla autoencoders.**

**Answer:** Vanilla autoencoders, also known as traditional or basic autoencoders, are a type of neural network architecture used for unsupervised learning and dimensionality reduction. They consist of an encoder network and a decoder network, which work together to learn a compressed representation of the input data.

Here's how vanilla autoencoders work:

1.  **Encoder**: The encoder network takes the input data and maps it to a lower-dimensional latent space representation. This is achieved by applying a series of transformations (typically linear and nonlinear) to the input data, gradually reducing its dimensionality until it reaches the desired size of the latent space.

2.  **Latent Space**: The latent space representation generated by the encoder serves as a compressed encoding of the input data. It captures the most important features or patterns present in the data while discarding redundant or less relevant information.

3.  **Decoder**: The decoder network takes the latent space representation produced by the encoder and maps it back to the original input space. Similar to the encoder, the decoder applies a series of transformations to the latent space representation to reconstruct the input data as accurately as possible.

4.  **Reconstruction Loss**: During training, vanilla autoencoders are optimized to minimize a reconstruction loss, which measures the difference between the input data and the reconstructed output. Common choices for the reconstruction loss include mean squared error (MSE) or binary cross-entropy, depending on the nature of the input data.

5.  **Training**: The autoencoder is trained on a dataset consisting of unlabeled examples of the input data. The encoder and decoder parameters are updated iteratively using backpropagation and gradient descent to minimize the reconstruction loss.

    Vanilla autoencoders are used for various tasks, including:

    *   Dimensionality reduction: By learning a lower-dimensional representation of the input data, autoencoders can be used for tasks such as data compression and visualization.
    *   Data denoising: Autoencoders can learn to reconstruct clean versions of noisy input data by training on pairs of noisy and clean examples.
    *   Anomaly detection: Autoencoders can learn to reconstruct normal patterns in the data and identify deviations from these patterns as anomalies.

While vanilla autoencoders are conceptually simple and easy to implement, they may struggle to learn complex or highly nonlinear transformations, especially when the input data distribution is complex. More advanced variants of autoencoders, such as convolutional autoencoders and variational autoencoders, address some of these limitations and offer improved performance for specific tasks.

2.  **What are Sparse autoencoders.**

**Answer:** Sparse autoencoders are a variant of autoencoder neural networks designed to learn sparse representations of input data. Unlike traditional autoencoders, which aim to learn dense representations

where most units in the hidden layers are active, sparse autoencoders encourage sparsity by penalizing the activation of hidden units, effectively forcing them to be inactive most of the time.

Here's how sparse autoencoders work:

1. **Sparse Constraints**: In sparse autoencoders, a sparsity constraint is imposed on the activations of the hidden units. This constraint encourages the hidden units to be mostly inactive, meaning they have low activation values most of the time.

2. **Regularization Term**: To enforce sparsity, a regularization term is added to the loss function of the autoencoder. This regularization term penalizes the deviation of the average activation of each hidden unit from a target sparsity value, often set to a small value close to zero.

3. **Training**: During training, the autoencoder learns to reconstruct the input data while simultaneously minimizing the reconstruction error and the sparsity regularization term. This encourages the model to learn a sparse representation that captures the most salient features of the input data.

   Sparse autoencoders have several advantages and applications:

   - **Dimensionality Reduction**: Like traditional autoencoders, sparse autoencoders can be used for dimensionality reduction, where the learned sparse representations encode the most important features of the input data in a compressed form.

   - **Feature Learning**: Sparse autoencoders can be used for unsupervised feature learning, where the sparse representations learned by the model can be transferred and fine-tuned for downstream supervised tasks, such as classification or regression.

   - **Anomaly Detection**: Sparse autoencoders can be effective for anomaly detection tasks, where the model learns to reconstruct normal patterns in the data and identifies deviations from these patterns as anomalies. Sparse representations make it easier to detect unusual or rare events.

Sparse autoencoders are particularly useful when the input data has high dimensionality or contains redundant information, as they can learn compact and informative representations that capture the essential features of the data. However, training sparse autoencoders can be more challenging and computationally expensive compared to traditional autoencoders due to the additional sparsity regularization term.

3. What are Denoising autoencoders.

**Answer:** Denoising autoencoders are a type of autoencoder neural network designed to learn robust representations of input data by reconstructing clean versions of noisy input data. They are trained to map corrupted or noisy input data to their clean counterparts, thereby learning to filter out noise and capture underlying patterns in the data.

Here's how denoising autoencoders work:

1. **Corrupted Input**: During training, denoising autoencoders are presented with input data that has been artificially corrupted by adding noise or introducing other types of perturbations. The corrupted input serves as the input to the autoencoder.

2. **Reconstruction Objective**: The autoencoder is trained to reconstruct the original, uncorrupted input data from the noisy input. The reconstruction objective is to minimize the difference between the clean input and the output reconstructed by the autoencoder.

3. **Noise Injection**: Different types of noise can be introduced to corrupt the input data, such as Gaussian noise, dropout noise, or masking noise. The choice of noise type depends on the characteristics of the data and the desired robustness of the learned representations.

4. **Training**: During training, the autoencoder learns to denoise the corrupted input by capturing the underlying structure and patterns in the data. By optimizing the reconstruction objective, the autoencoder learns to filter out the noise and produce clean reconstructions of the input data.

   Denoising autoencoders have several advantages and applications:

   - **Noise Robustness**: Denoising autoencoders learn to encode robust representations of input data that are less sensitive to noise and variations in the input. This makes them useful for tasks where input data may be corrupted or subject to noise, such as image denoising or speech recognition.

   - **Feature Learning**: By training on noisy input data, denoising autoencoders can learn features that are more invariant to irrelevant variations or nuisances in the data. This can lead to better generalization performance on downstream tasks, such as classification or regression.

   - **Anomaly Detection**: Denoising autoencoders can be used for anomaly detection, where they learn to reconstruct normal patterns in the data and identify deviations from these patterns as anomalies. Anomalies often manifest as distortions or inconsistencies in the reconstructed data.

Denoising autoencoders offer a principled approach to learning robust representations of input data by explicitly modeling and removing noise during training. They have been successfully applied in various domains, including computer vision, natural language processing, and signal processing, where data may be corrupted or subject to noise.

## 4. What are Convolutional autoencoders.

**Answer:** Convolutional autoencoders are a type of autoencoder neural network architecture that uses convolutional layers for both the encoder and decoder components. They are particularly well-suited for learning representations of image data due to their ability to capture spatial hierarchies of features.

Here's how convolutional autoencoders work:

1. **Encoder**: The encoder component of a convolutional autoencoder consists of convolutional layers followed by downsampling layers (e.g., max-pooling layers). These layers process the input image and progressively reduce its spatial dimensions while extracting hierarchical features.

2. **Latent Space**: The output of the encoder is a low-dimensional latent space representation of the input image. This representation captures the most salient features of the input image in a compact and compressed form.

3. **Decoder**: The decoder component of a convolutional autoencoder consists of convolutional layers followed by upsampling layers (e.g., transposed convolutional layers or nearest-neighbor upsampling). These layers gradually increase the spatial dimensions of the latent representation while reconstructing the original input image.

4. **Reconstruction Objective**: During training, the autoencoder is optimized to minimize the reconstruction error between the input image and the output reconstructed by the decoder. Common choices for the reconstruction loss include mean squared error (MSE) or binary cross-entropy, depending on the nature of the input data.

Convolutional autoencoders have several advantages and applications:

- **Image Denoising**: Convolutional autoencoders can be trained to denoise noisy images by reconstructing clean versions of the input images. The convolutional layers learn to capture spatial dependencies and patterns in the data, making them effective for filtering out noise.

- **Dimensionality Reduction**: Convolutional autoencoders can be used for dimensionality reduction tasks, where the learned latent space representation captures the most important features of the input images in a compressed form. This can be useful for tasks such as image compression or visualization.

- **Feature Learning**: Convolutional autoencoders can learn hierarchical representations of image data, where lower layers capture low-level features (e.g., edges, textures) and higher layers capture more abstract and complex features (e.g., objects, scenes). These learned features can be transferred and fine-tuned for downstream tasks such as image classification or object detection.

Convolutional autoencoders have been widely used in computer vision applications, including image denoising, image compression, image generation, and feature learning. They offer an effective and principled approach to learning compact and informative representations of image data.

## 5. What are Stacked autoencoders.

**Answer:** Stacked autoencoders, also known as deep autoencoders, are a type of autoencoder neural network architecture that consists of multiple layers of encoder and decoder units stacked on top of each other. They are used to learn hierarchical representations of input data by progressively transforming it into a lower-dimensional latent space and then reconstructing it back to the original input space.

Here's how stacked autoencoders work:

1. **Encoder**: Each layer of the stacked autoencoder consists of an encoder unit that maps the input data to a lower-dimensional representation. The encoder typically consists of multiple fully connected layers or convolutional layers followed by nonlinear activation functions. The output of each encoder layer serves as the input to the next encoder layer.

2. **Latent Space**: The output of the last encoder layer represents the low-dimensional latent space representation of the input data. This latent representation captures the most salient features of the input data in a compressed and abstract form.

3. **Decoder**: Each layer of the stacked autoencoder also consists of a decoder unit that maps the latent representation back to the original input space. The decoder mirrors the structure of the encoder, with fully connected layers or convolutional layers followed by nonlinear activation functions. The output of each decoder layer serves as the input to the next decoder layer.

4. **Reconstruction Objective**: During training, the stacked autoencoder is optimized to minimize the reconstruction error between the input data and the output reconstructed by the decoder. This is

typically done using a loss function such as mean squared error (MSE) or binary cross-entropy, depending on the nature of the input data.

5. **Training**: Stacked autoencoders are trained using unsupervised learning, where the input data is used to train the model without requiring labeled output data. The parameters of the encoder and decoder layers are learned iteratively using backpropagation and gradient descent to minimize the reconstruction error.

Stacked autoencoders have several advantages and applications:

- **Hierarchical Representation Learning**: By stacking multiple encoder and decoder layers, stacked autoencoders learn hierarchical representations of input data, where lower layers capture low-level features (e.g., edges, textures) and higher layers capture more abstract and complex features (e.g., objects, scenes).
- **Dimensionality Reduction**: Stacked autoencoders can be used for dimensionality reduction tasks, where the learned latent space representation captures the most important features of the input data in a compressed form. This can be useful for tasks such as data compression or visualization.
- **Unsupervised Pretraining**: Stacked autoencoders can be used for unsupervised pretraining of deep neural networks. Once trained on unlabeled data, the learned representations can be transferred and fine-tuned for downstream supervised tasks, such as classification or regression.

Stacked autoencoders have been successfully applied in various domains, including computer vision, natural language processing, and signal processing, where they have been used for tasks such as image reconstruction, feature learning, and anomaly detection.

6. Explain how to generate sentences using LSTM autoencoders.

**Answer:** Generating sentences using LSTM autoencoders involves training the autoencoder model to learn a compressed representation of input sentences and then using the decoder part of the model to generate new sentences based on the learned representations. Here's a step-by-step explanation of how to generate sentences using LSTM autoencoders:

1. **Data Preprocessing**:

   - Tokenize the input sentences into words or subword units (e.g., using the Tokenizer class in TensorFlow or Keras).
   - Convert the tokenized words into sequences of numerical indices (e.g., using the texts_to_sequences method).
   - Pad or truncate the sequences to a fixed length to ensure uniform input size (e.g., using the pad_sequences function).

2. **Architecture Design**:

   - Define an LSTM-based autoencoder architecture consisting of an encoder LSTM layer and a decoder LSTM layer.
   - The encoder LSTM layer processes the input sequences and produces a fixed-size latent representation.
   - The decoder LSTM layer takes the latent representation as input and generates sequences of words as output.

3. **Training**:

- Train the LSTM autoencoder using the input sequences as both input and target data.
- Use a reconstruction loss (e.g., mean squared error or cross-entropy) to measure the difference between the input sequences and their reconstructions.
- Optimize the model parameters using an optimization algorithm such as stochastic gradient descent (SGD) or Adam.

4. **Generating Sentences**:

- Given a seed input sequence or a randomly initialized latent representation, feed it into the decoder LSTM layer.
- Generate the next word in the sequence based on the output of the decoder LSTM layer.
- Repeat the process iteratively, feeding the generated words back into the decoder to produce longer sequences.
- Optionally, use techniques like temperature scaling or nucleus sampling to control the randomness of the generated text.

5. **Decoding Generated Sequences**:

- Convert the generated sequences of word indices back into words using the inverse mapping from step 1.
- Optionally, post-process the generated sentences to improve their readability or coherence (e.g., capitalization, punctuation).

6. **Evaluation and Fine-Tuning**:

- Evaluate the quality of the generated sentences using metrics such as BLEU score or human evaluation.
- Fine-tune the LSTM autoencoder architecture or hyperparameters based on the evaluation results to improve the quality of the generated text.

By following these steps, you can train an LSTM autoencoder to learn representations of input sentences and generate new sentences based on the learned representations. This approach is commonly used in natural language generation tasks such as text summarization, dialogue generation, and creative writing.

7. **Explain Extractive summarization.**

**Answer:** Extractive summarization is a technique used to create a summary of a longer text by selecting and extracting the most important sentences or passages directly from the original document. Rather than generating new sentences or paraphrasing the content, extractive summarization identifies and retains the most relevant information from the input text, preserving the original wording.

Here's how extractive summarization works:

1. **Text Preprocessing**: The input document is preprocessed to remove any unnecessary information, such as formatting, metadata, or boilerplate text. Tokenization may also be applied to break the text into individual sentences or passages.

2. **Sentence Scoring**: Each sentence or passage in the document is assigned a score that reflects its importance or relevance to the overall content. Various scoring methods can be used for this purpose, including:

- Word Frequency: Sentences containing frequently occurring words or phrases are considered more important.
- TF-IDF (Term Frequency-Inverse Document Frequency): Sentences with rare or distinctive terms that are important in the context of the document are prioritized.
- TextRank: A graph-based algorithm that represents sentences as nodes and computes their importance based on their connections to other sentences.

3. **Sentence Selection**: Based on the scores assigned to each sentence, a subset of sentences is selected to form the summary. The selection process can involve setting a threshold score and including all sentences above that threshold, or using a fixed number of top-scoring sentences.

4. **Summary Generation**: The selected sentences are concatenated to form the final summary. Optionally, post-processing steps such as sentence reordering or coherence improvement may be applied to enhance the readability and flow of the summary.

Extractive summarization has several advantages:

- **Preservation of Source Content**: Since extractive summarization directly selects and retains sentences from the original text, it ensures that the summary accurately represents the content and context of the source document.
- **Efficiency**: Extractive summarization is computationally efficient, as it does not involve generating new text or complex language modeling. It simply selects and outputs existing sentences from the input document.
- **Interpretability**: The summaries produced by extractive summarization are often more interpretable and easier to understand, as they maintain the original wording and structure of the source text.

However, extractive summarization also has limitations:

- **Redundancy**: Extractive summaries may contain redundant or repetitive information, especially if multiple sentences convey similar meaning or cover the same topic.
- **Incoherence**: Since sentences are selected independently based on their individual scores, the resulting summary may lack coherence or smooth transition between sentences.

8. **Explain Abstractive summarization.**

**Answer:** Abstractive summarization is a technique used to generate a concise summary of a longer text by interpreting and paraphrasing the original content in a new way. Unlike extractive summarization, which selects and retains existing sentences from the input document, abstractive summarization involves synthesizing new sentences that capture the essential information and meaning of the source text.

Here's how abstractive summarization works:

1. **Understanding the Input Text**: The first step in abstractive summarization is to understand the content and context of the input document. This may involve using natural language processing (NLP) techniques to parse and analyze the text, identifying key concepts, entities, and relationships.

2. **Generating a Summary**: Once the input text has been processed and understood, the abstractive summarization model generates a summary by:

- Generating new sentences that convey the main ideas and key points of the original text.

- Paraphrasing and rephrasing the content in a more concise and coherent manner.
- Ensuring that the summary maintains the overall meaning and tone of the original text while using different wording and structure.

3. **Language Generation**: Abstractive summarization models often employ advanced language generation techniques, such as recurrent neural networks (RNNs), long short-term memory (LSTM) networks, or transformer models (e.g., GPT, BERT). These models learn to generate text by predicting the most probable next word or sequence of words given the context provided by the input text.

4. **Training**: Abstractive summarization models are trained on pairs of input documents and their corresponding human-written summaries. During training, the model learns to generate summaries that are similar to the human-written summaries, optimizing for metrics such as semantic similarity, fluency, and coherence.

5. **Evaluation and Fine-Tuning**: The generated summaries are evaluated using metrics such as ROUGE (Recall-Oriented Understudy for Gisting Evaluation), BLEU (Bilingual Evaluation Understudy), or human judgment. The model parameters may be fine-tuned based on the evaluation results to improve the quality of the generated summaries.

Abstractive summarization has several advantages:

- **Flexibility**: Abstractive summarization allows for more flexibility in summarizing the content, as it can generate new sentences and rephrase the original text to produce summaries that are more concise and coherent.
- **Reduced Redundancy**: Unlike extractive summarization, which may include redundant or repetitive information, abstractive summarization can generate summaries that contain only the most relevant and important points.
- **Improved Readability**: Abstractive summaries are often more readable and easier to understand, as they are crafted to convey the main ideas of the original text in a more natural and fluent manner.

However, abstractive summarization also has challenges:

- **Complexity**: Generating high-quality abstractive summaries requires sophisticated language understanding and generation capabilities, which can be challenging to achieve, especially for longer and more complex documents.
- **Semantic Accuracy**: Ensuring that the generated summaries accurately capture the meaning and context of the original text can be difficult, particularly when dealing with ambiguous or nuanced language.

## 9. Explain Beam search.

**Answer:** Beam search is a heuristic search algorithm used in natural language processing (NLP) and other fields to find the most likely sequence of outputs from a probabilistic model, such as a language model or a machine translation system. It is commonly used in tasks like text generation, where the goal is to produce fluent and coherent sequences of words.

Here's how beam search works:

1. **Initialization**: Beam search begins with an initial set of candidate sequences, often referred to as the beam. This set typically contains only one sequence initially, which is usually the start-of-sequence token (e.g., **\<start\>** or **\<s\>**).

2. **Expansion**: At each step of the search, the beam is expanded by generating the next set of candidate sequences based on the probabilities assigned to each possible output token by the probabilistic model. For each candidate sequence in the current beam, the model generates a set of possible next tokens along with their associated probabilities.

3. **Selection**: The expanded set of candidate sequences is sorted based on their probabilities, and only the top-k sequences with the highest probabilities are retained, where k is a hyperparameter known as the beam width. This step ensures that the search space is limited to the most promising sequences.

4. **Repeat**: Steps 2 and 3 are repeated iteratively until a stopping criterion is met, such as generating an end-of-sequence token (e.g., **\<end\>** or **\<\s\>**), reaching a maximum length, or reaching a maximum number of iterations.

5. **Termination**: Once the stopping criterion is met, the beam search algorithm selects the sequence with the highest probability from the final beam as the output sequence.

   Beam search allows the model to explore multiple possible sequences simultaneously, taking into account the probabilities assigned to each token at each step. By retaining only the most promising sequences at each step, beam search can efficiently find high-quality output sequences while avoiding exhaustive search through the entire space of possible sequences.

   However, beam search has some limitations:

   - **Suboptimal Solutions**: Beam search may produce suboptimal solutions, especially when the true optimal sequence diverges from the most likely sequence at early steps of the search.
   - **Exponential Growth**: The number of candidate sequences grows exponentially with the beam width, leading to increased computational complexity and memory requirements, especially for large beam widths.

## 10. Explain Length normalization.

**Answer:** Length normalization is a technique used to address the bias introduced by beam search towards shorter output sequences, particularly in tasks like text generation, where the length of the generated sequences varies. In beam search, shorter sequences tend to have higher probabilities because they have fewer steps to accumulate probabilities, leading to a bias towards generating shorter outputs.

Length normalization aims to mitigate this bias by adjusting the probabilities of candidate sequences based on their lengths. The idea is to penalize shorter sequences and reward longer sequences, so that the probabilities of all sequences are more comparable regardless of their lengths.

Here's how length normalization works:

1. **Compute Length Penalties**: For each candidate sequence in the beam, compute a length penalty that adjusts its probability based on its length. There are several ways to define the length penalty, but a common approach is to use a simple length penalty function, such as:

   length penalty=$(c+$length of sequence$c+1)\alpha$length penalty=$(c+1c+$length of sequence$)\alpha$

where $c$ is a hyperparameter that controls the strength of the penalty, and $\alpha$ is another hyperparameter that controls the steepness of the penalty curve. Larger values of $c$ result in stronger penalties for shorter sequences, while larger values of $\alpha$ result in steeper penalty curves.

2. **Apply Length Penalties**: Multiply the probabilities of candidate sequences by their corresponding length penalties. This adjusts the probabilities of shorter sequences downward and the probabilities of longer sequences upward, making them more comparable.

3. **Rescale Probabilities**: After applying length penalties, rescale the probabilities of all candidate sequences so that they sum to one. This ensures that the adjusted probabilities remain valid probability distributions.

   By applying length normalization, beam search can produce output sequences that are more balanced in terms of their lengths and less biased towards shorter sequences. This can lead to more diverse and coherent outputs, especially in tasks where generating outputs of varying lengths is desirable, such as text generation and machine translation.

Length normalization is commonly used as a post-processing step in conjunction with beam search and other decoding algorithms in natural language processing tasks. It helps mitigate the bias towards shorter sequences and improves the overall quality and diversity of the generated outputs.

## 11. Explain Coverage normalization.

**Answer:** Coverage normalization is a technique used in neural machine translation (NMT) to address the problem of repeated translations or over-translation of certain source words or phrases. It aims to encourage the model to generate translations that cover all the information present in the source sentence, thus improving the overall coherence and fluency of the translated output.

In NMT, attention mechanisms are often used to align source and target words during translation. Coverage normalization extends the attention mechanism to keep track of the attention weights assigned to each source word over multiple decoding steps and penalize repeated attention to the same source words.

Here's how coverage normalization works:

1. **Initialization**: At each decoding step, a coverage vector $c_t$ is initialized to zero. This vector has the same length as the source sentence and keeps track of the attention coverage of each source word.

2. **Attention Calculation**: The attention mechanism computes the attention weights $\alpha_t$ for each source word at the current decoding step $t$ based on the current decoder state $s_t$ and the source context vectors $h_i$:

   $$\alpha_t(i) = \text{softmax}(v_a^T \tanh(W_a[s_t; h_i] + b_a))$$

3. **Coverage Update**: The coverage vector $c_t$ is updated based on the attention weights $\alpha_t$ computed at the current step:

   $$c_t = c_{t-1} + \alpha_t$$

   This update accumulates the attention weights over time, indicating how much each source word has been attended to up to the current decoding step.

4. **Coverage Penalty**: To encourage coverage of all source words, a coverage penalty term is added to the loss function during training. This penalty term measures the discrepancy between the current coverage vector $c_t$ and a uniform distribution, indicating equal coverage of all source words. A common choice for the coverage penalty is the Kullback-Leibler (KL) divergence:

$$\text{Coverage penalty}=\sum_i \text{KL}(\text{uniform}(0,1)\|c_t(i))$$

This penalty encourages the model to distribute its attention more evenly across all source words and discourage repeated attention to the same words.

By incorporating coverage normalization into the attention mechanism, NMT models can generate more diverse and coherent translations by avoiding over-translation of certain source words or phrases. This leads to improvements in translation quality, especially for long or complex sentences where coverage issues are more likely to occur.

## 12. Explain ROUGE metric evaluation.

**Answer:** ROUGE (Recall-Oriented Understudy for Gisting Evaluation) is a set of metrics used to evaluate the quality of summaries produced by automatic summarization systems. It measures the overlap between the generated summary and one or more reference summaries (often written by humans) in terms of shared n-grams, word sequences, or semantic units. ROUGE is widely used in natural language processing (NLP) and information retrieval (IR) research to assess the effectiveness of summarization algorithms.

ROUGE consists of several evaluation metrics, each focusing on a different aspect of summary quality:

1. **ROUGE-N**: Measures the overlap of n-grams (contiguous sequences of n words) between the generated summary and the reference summaries. It typically computes precision, recall, and F1-score for n-grams of various lengths (unigrams, bigrams, trigrams, etc.).

2. **ROUGE-L**: Measures the longest common subsequence (LCS) between the generated summary and the reference summaries. It considers not only individual words but also their ordering and sequence structure.

3. **ROUGE-W**: Similar to ROUGE-L, but allows for skipping words in the LCS to account for word insertions, deletions, or substitutions. This metric penalizes more heavily for longer gaps between matching words.

4. **ROUGE-S**: Measures the skip-bigram overlap between the generated summary and the reference summaries. Skip-bigrams are pairs of words that occur in the summary within a certain window size, regardless of their relative positions.

5. **ROUGE-SU**: An extension of ROUGE-S that considers skip-bigrams with varying gap sizes. This metric captures a wider range of semantic similarity between the summary and the reference summaries.

For each ROUGE metric, precision, recall, and F1-score are typically computed to evaluate the quality of the generated summary. Precision measures the proportion of generated n-grams (or other units) that appear in the reference summaries, recall measures the proportion of reference n-grams that appear in the generated summary, and F1-score is the harmonic mean of precision and recall.

ROUGE evaluation involves comparing the generated summary against one or more reference summaries using the chosen ROUGE metrics. Higher ROUGE scores indicate better agreement between the generated summary and the reference summaries, suggesting higher quality summarization. ROUGE scores are often

reported as averages over multiple documents or datasets to provide a comprehensive evaluation of summarization system performance.

Overall, ROUGE metrics provide a quantitative and standardized way to evaluate the effectiveness of automatic summarization systems, helping researchers and practitioners assess and compare different algorithms and techniques for summarization.