

Article

Not peer-reviewed version

---

# Semantic Information Theory: Recent Advances and Future Challenges

---

[Gangtao Xin](#) , [Pingyi Fan](#) <sup>\*</sup> , Khaled B. Letaief

Posted Date: 19 October 2023

doi: 10.20944/preprints202310.1208.v1

Keywords: Semantic information theory; semantic communication; semantic distortion; 6G; goal-oriented communications; joint source-channel coding; deep learning; information bottleneck






Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

## Article

# Semantic Information Theory: Recent Advances and Future Challenges

Gangtao Xin <sup>1,2</sup> , Pingyi Fan <sup>1,2,\*</sup>  and Khaled B. Letaief <sup>3</sup> 

<sup>1</sup> Department of Electronic Engineering, Tsinghua University, Beijing 100084, China; xgt19@mails.tsinghua.edu.cn

<sup>2</sup> Beijing National Research Center for Information Science and Technology, Tsinghua University, Beijing 100084, China

<sup>3</sup> Department of Electrical and Computer Engineering, Hong Kong University of Science and Technology (HKUST), Hong Kong

\* Correspondence: fpy@tsinghua.edu.cn; Tel.: +86-010-6279-6973

**Abstract:** In recent years, semantic communication has received significant attention from both academia and industry, driven by the growing demands for ultra-low latency and high-throughput capabilities in emerging intelligent services. Nonetheless, a comprehensive and effective theoretical framework for semantic communication has yet to be established. In particular, finding the fundamental limits of semantic communication, exploring the capabilities of semantic-aware networks, or utilizing theoretical guidance for deep learning in semantic communication are very important but yet still unresolved issues. In this paper, we delve into the pertinent advancements in semantic information theory. Grounded in the foundational work of Claude Shannon, we present the latest developments in semantic entropy, semantic rate distortion, and semantic channel capacity. Additionally, we will analyze some open problems in semantic information measurement and semantic coding, providing a theoretical basis for the design of a semantic communication system. Furthermore, we carefully review several mathematical theories and tools and evaluate their applicability in the context of semantic communication. Finally, we shed light on the challenges encountered in both semantic communication and semantic information theory.

**Keywords:** semantic information theory; semantic communication; semantic distortion; 6G; goal-oriented communications; joint source-channel coding; deep learning; information bottleneck

## 1. Introduction

In recent years, the rapid development of wireless communications and the increasing demand for intelligent processing have given rise to remarkable growth in various emerging intelligent services. However, this surge has brought new challenges to communication and computing technology. On the one hand, the success of these emerging intelligent businesses, such as the industrial internet, virtual/augmented/mixed reality, metaverse, and holographic communications, heavily relies on training large foundational models with extensive datasets. The substantial traffic generated by these new applications has the potential to overload existing communication networks. Consequently, it is imperative that the communication infrastructure incorporates intelligence to efficiently handle this traffic in a timely and organized manner. On the other hand, these intelligent services require extremely low end-to-end latency. For instance, in the realm of autonomous driving, vehicles depend on near-instantaneous data exchange to make split-second decisions, thereby avoiding potential traffic accidents. Similarly, in the context of remote surgery systems, the timely update of surgical tool positions is necessary to ensure the safety and precision of medical procedures. As a result, communication technology must take into account the relevance and urgency of traffic, enabling the swift and reliable extraction and delivery of task-related information. This paradigm shift highlights the importance of the evolution of communication network architecture, moving beyond a sole focus on high-speed symbol transmission and towards prioritizing high-quality semantic exchange [1–3].

Semantic communication is a novel architecture that seamlessly integrates tasks and user requirements into the communication process, which is expected to greatly improve communication efficiency and user experience. It emphasizes the efficient exchange of semantics and the clear communication of meaning. Furthermore, this innovative paradigm has the potential to fundamentally address the complex compatibility issues that have plagued traditional data-based communication systems, including challenges spanning cross-system, cross-protocol, and cross-network domains. In the pioneering work of Weaver in 1953 [4], it was articulated that general communications involve problems at three levels, outlined as follows:

- Level A: How accurately can the symbols of communication be transmitted? (The technical problem.)
- Level B: How precisely do the transmitted symbols convey the desired meaning? (The semantic problem.)
- Level C: How effectively does the received meaning affect conduct in the desired way? (The effectiveness problem.)

Figure 1 illustrates a visual representation of the three-level communication architecture and its underlying mechanism. Moving from the bottom to the top, one encounters the technical level, semantic level, and effectiveness level of communication. From the technical layer to the semantic layer, the goal of communication shifts from the accurate transmission of data to the effective exchange of semantics embedded in data. These evolving communication goals necessitate corresponding changes in mathematical theory. The classic Shannon’s information theory, rooted in probability and statistics, primarily addresses the technical layer’s concerns, such as data compression and communication transmission rate. However, it falls short when applied to the semantic layer since it disregards the semantics of data information and fails to account for crucial semantic factors such as task relevance, time constraints, and urgency. In general, the mathematical theory of semantic communication and the mathematical representation of data semantics can be attributed to the problem of *semantic information theory*. While a recognized and unified theoretical framework for semantic information theory is currently absent, both academia and industry have experienced a significant increase in research activities in this domain in recent years. These endeavors have generated a need for the systematic organization and summarization of existing research findings, which can serve as a catalyst for further exploration and advancement in the field of semantic information theory.

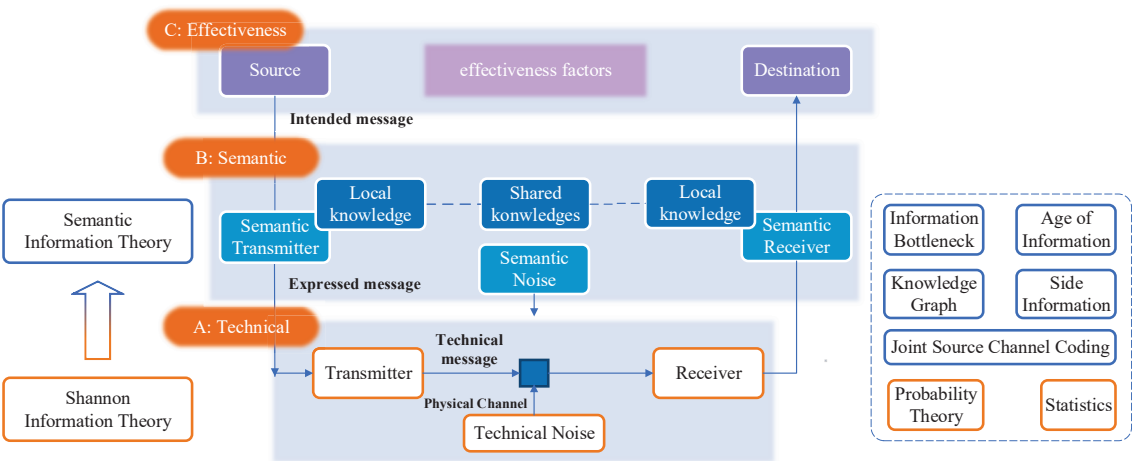


Figure 1. A three-level communication architecture.

While several reviews have delved into the realm of semantic communication. In [5–13], they have primarily centered on aspects related to systems, algorithms, and architectures, as well as their connections with deep learning. This article takes a distinctive perspective by focusing on the theoretical dimension—specifically, semantic information theory. From this point, we aim to

comprehensively review and examine recent advancements and chart the future directions of semantic information theory. Grounded in the foundational work of Claude Shannon, we present the latest developments in semantic entropy, semantic rate distortion, and semantic channel capacity. Moreover, we establish connections between semantic information theory and Shannon's information theory, with a primary focus on some core concepts of semantic information theory. Furthermore, we introduce various mathematical theories and tools, including concepts like the Age of Information (AoI) and Information Bottleneck (IB), which hold significant potential to propel semantic information theory forward.

The rest of this article is structured as follows. Section 2 provides an introduction to semantics and semantic communication. Sections 3 to 5 constitute an in-depth exploration of the fundamental concepts within semantic information theory. These sections cover essential topics, including semantic entropy, semantic rate distortion, and semantic channel capacity. In Section 6, we delve into the mathematical theories and tools that are relevant to the domain of semantic communication. Section 7 discusses the potential challenges that may arise during the development of semantic communication and semantic information theory. Finally, we conclude the paper in Section 8.

## 2. Semantic Communication

Although this paper serves as a summary of the latest research advancements in semantic information theory, it is necessary to establish an intuitive understanding of semantics and semantic communication before delving into the theoretical part. In this section, we present a brief introduction to semantic communication and introduce a general semantic communication system.

### 2.1. What is Semantic Communication?

Semantic communication serves distinct motivations and purposes compared to traditional digital communication. In the landmark of Shannon's 1948 paper, Shannon [14] stated that

*The fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point.*

While Weaver [4] emphasized that

*The semantic problems are concerned with the interpretation of meaning by the receiver, as compared with the intended meaning of the sender.*

These two statements correspond to different levels of communication. Shannon's sentence addresses the technical layer, specifically digital communication, while Weaver's sentence focuses on the semantic layer, namely, semantic communication. By comparing these two statements, we can find that the objective of semantic communication is not to replicate the transmitted messages, whether exact or approximate, but rather to convey their interpretations accurately. For example, consider the following conversation:

Alice: "Do you like bananas?"

Bob: "No, I hate eating any fruit."

In this conversation, Alice serves as a semantic sender while Bob assumes the role of a semantic receiver. Bob is able to interpret the meanings of the received message and relate it to his existing vocabulary. He knows that "hate" is an antonym of "like", and "banana" falls under the category of "fruit". Consequently, he can infer that "hate eating any fruit" implies "do not like bananas", despite the fact that the two statements have distinct syntactical structures [15]. Now consider a conversation between three persons:

Alice: "Bob, does Carol like bananas?"

Bob: "Carol, if you enjoy bananas?"

Carol: "No, I do not enjoy any fruit."

In this context, Bob acts as a semantic channel between Alice and Carol. While Bob may not precisely replicate Alice's original message, he adeptly retains its intended meaning. While assessing

the success of communication in a purely literal sense, there might be an engineering failure. However, there is no failure at the semantic level.

From these two examples, we can see that the objective of semantic communication lies in the effective exchange of semantics. In other words, whether the meaning carried by the symbol can be understood by the receiver. This model of communication capitalizes on the participants’ perception of the world and their responses to various concepts, thereby giving symbols a deeper and more abstract connotation. In summary, semantic communication is not to reproduce, but to convey after understanding.

2.2. What is a Semantic Communication System?

In general, current semantic communication systems are constructed upon digital communication frameworks. In other words, these systems still depend on the physical channel to transmit specific symbols and are not entirely detached from Shannon’s paradigm, which is consistent with **Weave’s** viewpoint. A semantic communication system usually comprises the following key components [7]:

- Semantic Encoder: This component is responsible for detecting and extracting the semantic content from the source message. It may also compress or eliminate irrelevant information to enhance efficiency.
- Channel Encoder: The role of the channel encoder is to encode and modulate the semantic features of the message as signal to combat any noise or interference that may occur during transmission.
- Channel Decoder: Upon receiving the signal, the channel decoder demodulates and decodes it, recovering the transmitted semantic features.
- Semantic Decoder: The semantic decoder interprets the information sent by the source and converts the received signal features into a format that is comprehensible to the destination user.
- Knowledge Base: The knowledge base serves as a foundation for the semantic encoder and decoder, enabling them to understand and infer semantic information accurately and effectively.

In general, a semantic communication system includes the five mentioned components, with the flexibility to add more as required for specific tasks. Figure 2 illustrates a general semantic communication architecture for the transmission task of image recognition. Rather than transmitting bit sequences that represent the entire image, the semantic transmitter in this architecture extracts only the features crucial for recognizing the object—in this case, a dog—from the source. Irrelevant information, like the image background, is intentionally omitted to minimize the transmitted data while maintaining performance quality [5]. In this model, the knowledge base empowers the semantic encoder and decoder to generate and reconstruct semantics related to image recognition, respectively.

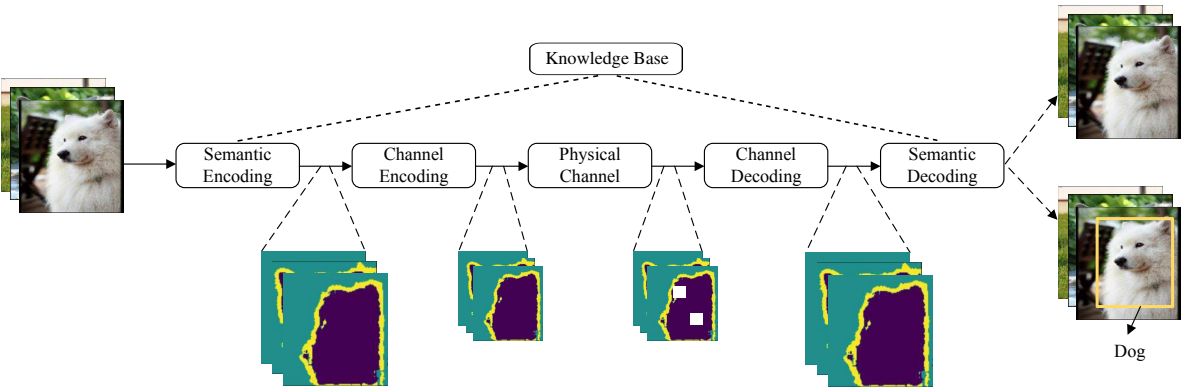


Figure 2. A general semantic communication architecture [5]

In summary, a semantic communication system is not a hypothetical concept but is built on a digital communication system, as it relies on physical channels to transmit essential symbols. The



information carried by symbols results from the empowerment of the knowledge base. In the following section, we will introduce several important concepts and theorems within semantic information theory, while also highlighting its distinctions from classical information theory.

### 3. Semantic Entropy

Entropy, which measures the uncertainty of a random variable, constitutes a fundamental concept in Shannon's information theory. Likewise, the quantification of semantic information forms the cornerstone of semantic information theory, referred to as *semantic entropy*. Semantic entropy serves as a metric for quantifying semantic uncertainty or the amount of semantic information. However, formulating an intuitive and universally applicable mathematical description of semantic entropy remains a formidable task. On the one hand, the semantic connotation is elusive to define and quantify. On the other hand, the generation mechanisms and processes of semantics remain obscure [3,8,16]. In this section, we delve into the essence of semantics and examine various definitions of semantic entropy.

#### 3.1. Statistical and Logical Probability

Let  $X$  be a discrete random variable with alphabet  $\mathcal{X}$  and probability mass function  $p(x) = \Pr\{X = x\}, x \in \mathcal{X}$ . In Shannon's information theory, the *entropy*  $H(X)$  of a discrete random variable  $X$  is defined by

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x). \quad (1)$$

In Equation (1),  $p(x)$  is the statistical probability of  $x$ , reflecting its frequency information. However, statistical probability is no longer the exclusive mathematical tool of choice for semantic communication. In general, probability includes two aspects, one is about logical probability and the other is about statistical probability [17].

- **Logical Probability:** Logical probability pertains to the degree of confirmation of a hypothesis with respect to an evidence statement, such as an observation report. A sentence regarding this concept relies on logical analysis rather than the direct observation of facts.
- **Statistical Probability:** Statistical probability refers to the relative frequency (in the long run) of one property of events or things with respect to another. A sentence concerning this concept is grounded in factual and empirical observations.

Shannon chose statistical probability as the mathematical foundation of information theory due to its ability to leverage principles from probability theory and stochastic processes. Through the application of the law of large numbers, Shannon can derive asymptotic equipartition property, paving the way for the derivation of several key theorems in information theory. Nonetheless, when it comes to semantic entropy, there is no widely accepted consensus on whether to employ statistical probability or logical probability. In the following discussion, we will see that semantic entropy was initially formulated using logical probability and subsequently evolved into various distinct formulations.

#### 3.2. Semantic Entropy

Semantic entropy originates from the analysis of natural language. In 1952, Carnap and Bar-Hillel [18] proposed the concept of semantic entropy as a means to quantify the amount of semantic information conveyed by a sentence. This concept aimed to assess the depth of meaning within a sentence, capturing its richness in conveying information. Let  $m(e)$  be the logical probability of event  $e$ , which signifies the likelihood that the sentence holds true in all possible situations. Then, the semantic entropy  $H_s(e)$  is defined by

$$H_s(e) = - \log m(e). \quad (2)$$

It is evident that the higher the logical probability of a sentence, the lower the semantic entropy. However, this gives rise to a paradox. Any statement that contradicts itself will possess an infinite amount of semantic information, such as  $H_s(e, \neg e)$  (" $\neg e$ " represents the counter-event or complementary event of " $e$ "), which becomes infinite.

In 2004, Floridi [19] proposed the strong semantic information theory, which resolved this paradox by utilizing the distance from actual events to represent the quantity of information. In 2011, D'Alfonso [20] provided a quantitative description of semantic information based on truthlikeness. Both Floridi and D'Alfonso measured the semantic information of a particular event relative to a reference event, yielding a value ranging from 0 to 1. However, these measurements heavily depend on the existence of reference events. Without a reference event, it becomes impossible for them to quantify the semantic entropy. Essentially, their work provides a measure of semantic similarity between two sentences, rather than a gauge of semantic uncertainty or informativeness. In alignment with Carnap's definition, several works have enriched and provided a more specific representation of semantic entropy by extending the connotation of  $m(e)$  [21]. In 2011, Bao [15] used propositional logic to expand the representation of  $m(e)$ . For a message (sentence)  $e$ , let  $W_e$  be the set of its models, i.e., worlds in which  $x$  is "true",  $W_x = \{w \in W | w \vdash x\}$  ( $\vdash$  is the logical entailment relation). Let  $\mu(w)$  be the statistical probability of model  $w$ . Then, the logical probability of  $e$  is

$$m(e) = \frac{\mu(W_e)}{\mu(W)} = \frac{\sum_{w \in W, w \vdash e} \mu(w)}{\sum_{w \in W} \mu(w)}. \quad (3)$$

Equation (3) shows that what matters is the total probability of the sentence model, not the model set's cardinality. In addition to the aforementioned research grounded in logical probability, we classify the work on semantic entropy into several distinct categories.

1) **Task-oriented:** The meaning and mechanism of semantic entropy should have various representations to suit different tasks. Chattopadhyay et al. [22] proposed the quantification of task-related semantic entropy, defined as the minimum number of semantic queries about data  $X$  required to solve the task  $V$ . It can be expressed as

$$H_s(X; V) \triangleq \min_E \mathbb{E}[|\text{Code}_Q^E(X)|], \quad (4)$$

where  $\text{Code}_Q^E(x)$  represents the query vector extracted from  $x$  using semantic encoder  $E$ .

For translation tasks, Melamed [23] proposed a method to measure the semantic entropy of words in a text. Specifically, let  $w$  represent a given word, the semantic entropy can be expressed as

$$H_s(w) = H(T|w) + N(w) = \sum_{t \in T} p(t|w) \log p(t|w) + p(\text{Null}|w) \log F(w). \quad (5)$$

Among these components,  $H(T|w)$  represents translation inconsistency, signifying the uncertainty when translating a word, where  $T$  represents the set of target words.  $N(w)$  reflects the impact of empty links for word  $w$ , indicating the likelihood of encountering translation difficulties between languages.  $F(w)$  represents the frequency of word  $w$ , and  $p(\text{Null}|w)$  is the probability of encountering problems when translating  $w$ .

For classification tasks, Liu et al. [24] introduced the concepts of matching degree and membership degree to define semantic entropy. Membership degree, a concept from fuzzy set theory [25], is challenging to express analytically. It is generally given based on experience. If we denote  $\mu_\zeta(x)$  as the membership degree of each  $x \in X$ , then for a specific category  $C_j$ , the matching degree is defined as

$$D_j(\zeta) = \frac{\sum_{x \in X_{C_j}} \mu_\zeta(x)}{\sum_{x \in X} \mu_\zeta(x)} \quad (6)$$

For category  $C_j$ , its semantic entropy is defined as  $H_{C_j}(\zeta) = -D_j(\zeta) \log D_j(\zeta)$ . To obtain the overall semantic entropy for  $X$ , one can sum the semantic entropy contributions from all categories, which is expressed as

$$H_s(\zeta) = \sum_j H_{C_j}(\zeta) \quad (7)$$

2) **Knowledge-based:** Semantics involves the comprehension of symbols, and knowledge plays a crucial role in the process of semantic encoding and representation. Choi et al. [26] explored the semantic entropy of a sentence from the perspective of knowledge bases using logical probability. Let the knowledge base be denoted as  $K$ . Let  $m[K \vdash e]$  be the probability that  $e$  is true relative to the knowledge base  $K$ , which can be simplified as  $m_e$ . Then the semantic entropy of  $e$  relative to  $K$  is calculated as:

$$H_s(e) = -m_e \log m_e + (1 - m_e) \log(1 - m_e). \quad (8)$$

This equation quantifies the semantic entropy of  $e$  with respect to the knowledge base  $K$ .

Moreover, expansion has the capability to amalgamate simple elements into complex systems, potentially leading to the emergence of intelligence. In the realm of human language, sentences are constructed from components such as subjects, predicates, objects, and attributive complements, enabling the expression of profound meanings that single words cannot convey. Xin and Fan [27] advocated for the extensibility of semantics, emphasizing that the representation of semantic entropy should encompass the notion of expansion. As semantics expand, knowledge often collisions. This can be likened to the phenomenon where as a country expands its territory, armed conflicts may arise. Semantics is a product born from the interaction between knowledge and signals. For instance, while "Apple Inc." falls under the category of a business company, and "fifteenth" is a numerical concept, their collision can give rise to a new word - "iPhone", signifying a mobile communication product. Let  $X_1$  and  $X_2$  denote signals, and let  $K_A^1$  and  $K_A^2$  represent two instances of knowledge. Let  $T$  and  $\hat{T}$  be the semantics of transmitter and receiver, respectively. Then, one step of the expansion architecture of semantic communication is described:

$$\begin{array}{ccc} X & \xrightarrow{(a)} & \hat{X} \\ \uparrow & & \downarrow \\ T \leftarrow X_1 \oplus X_2 & & \hat{X}_1 \oplus \hat{X}_2 \rightarrow \hat{T} \\ \uparrow (c) & & \uparrow (d) \\ K_A^1 \odot K_A^2 & \xrightarrow{(b)} & K_B^1 \odot K_B^2 \end{array} \quad (9)$$

where (a) is the explicit channel, (b) is the implicit channel, and (c) and (d) are semantic encoding and decoding, respectively. The semantic entropy can be expressed as  $H_s(X_1 \oplus X_2, K_A^1 \otimes K_A^2)$ , where  $\oplus$  denotes expansion and  $\otimes$  represents collision.

3) **Context-related:** Different from logical probability or statistical probability models, the forms of derivation for semantic entropy vary depending on the specific context. Kowlchinsky and Wolpert [28] defined semantic information as grammatical information that describes the relationship between a system and its environment. Kountouris and Pappas [8] used Rényi entropy [29] to measure semantic information. Venhuizen et al. [30] derived semantic entropy from a language understanding model grounded in background knowledge. Lu [31] introduced general information theory and employed concepts such as the Bayesian formula, logical probability, and fuzzy sets to mathematically describe semantic information.

#### 4. Semantic Rate Distortion

In the communication process, achieving a perfect performance is not always possible. It is conceivable for the receiver to obtain symbols that do not align with those sent by the sender. Additionally, describing an arbitrary real number necessitates an infinite number of bits. Therefore,



representing a continuous random variable with a finite representation can never be flawless. To approach this question appropriately, it is necessary to establish the quality of a source's representation. This is achieved by defining a distortion measure, which serves as a metric for evaluating the disparity between the random variable and its representation. In Shannon's information theory, rate-distortion theory addresses the coding problem of continuous random variables or vectors in the presence of distortion [32].

If we possess a source capable of generating a sequence  $X_1, X_2, \dots, X_n$ , i.i.d.  $\sim p(x), x \in \mathcal{X}$ . The encoder describes the source sequence  $X^n$  by an index  $f_n(X^n) \in \{1, 2, \dots, 2^{nR}\}$ . The decoder represents  $X^n$  by an estimate  $\hat{X}^n \in \hat{\mathcal{X}}^n$ . A distortion function, or distortion measure, is a mathematical mapping denoted as  $d: \mathcal{X} \times \hat{\mathcal{X}} \rightarrow \mathcal{R}^+$ . It operates on pairs of elements, each drawn from the source alphabet  $\mathcal{X}$  and the reproduction alphabet  $\hat{\mathcal{X}}$ , and produces non-negative real numbers as its output. The distortion, denoted as  $d(x, \hat{x})$ , quantifies the cost associated with representing the symbol  $x$  using the symbol  $\hat{x}$ . When considering sequences, such as  $x^n$  and  $\hat{x}^n$ , the distortion  $d(x^n, \hat{x}^n)$  is calculated as the average per-symbol distortion across the elements of the sequence. Based on this, one can define the *rate distortion function*.

**Definition 1.** The rate distortion function for a source  $X \sim p(x)$  and distortion measure  $d(x, \hat{x})$  is

$$R(D) = \min_{p(\hat{x}|x)} I(X; \hat{X}), \quad (10)$$

where the minimization is over all conditional contributions  $p(\hat{x}|x)$  for which the joint distribution  $p(x, \hat{x}) = p(x)p(\hat{x}|x)$  satisfies the expected distortion constraint  $\mathbb{E}(d(X, \hat{X})) \leq D$ .

The value  $R(D)$  represents the infimum of rates  $R$  achievable for a given distortion  $D$ . Building upon the foundation of the rate distortion function, Shannon subsequently derived the influential rate distortion theorem.

**Theorem 1.** (Rate distortion theorem) If  $R > R(D)$ , there exists a sequence of codes  $\hat{X}^n(X^n)$  with the number of codewords  $|\hat{X}^n(\cdot)| \leq 2^{nR}$  and  $\mathbb{E} d(X^n, \hat{X}^n(X^n)) \rightarrow D$ . If  $R \leq R(D)$ , no such codes exist.

The rate distortion theorem addresses two fundamental results regarding distortions: Firstly, given a source distribution and a distortion measure, it determines the minimum expected distortion achievable at a specific rate. Secondly, it establishes the minimum rate description required to achieve a particular distortion. In this section, we delve into the semantic level of rate distortion theory, drawing upon Shannon's foundation work. We explore topics such as semantic mismatch, the semantic rate distortion theorem, and semantic coding.

In the field of semantic information theory, the exploration of semantic rate distortion holds significant importance. In related studies, two distortion measurement functions, namely, semantic distortion and symbol distortion, are commonly employed to assess the effects of coding on communication quality. Utilizing established semantic evaluation criteria, defining and implementing semantic rate distortion becomes more accessible when compared to the complexities of semantic entropy and semantic channel capacity. Next, we first introduce the semantic mismatch evaluation criteria for various objects.

#### 4.1. Metrics for Semantic Mismatch

In general, the quality of semantic communication diverges from the traditional measure of bit error rate (BER) and symbol error rate (SER) commonly used in digital communication. Semantic communication often employs metrics capable of assessing semantic similarity, which aligns more closely with human perception. Furthermore, there is no universal metric for semantic mismatch. One-size-fits-all is unrealistic and impossible in semantic communication. It generally adapts to specific

tasks for information sources. In this work, we will concentrate on a select set of representative metrics. We begin by introducing the performance metrics employed in contemporary semantic communication systems for images, text, and audio, respectively.

1) **Image:** The measurement of similarity between two images, denoted as  $A$  and  $B$ , is expressed as follows:

$$\mathcal{L}(A, B) = \|f(A) - f(B)\|_2^2, \quad (11)$$

where  $f(\cdot)$  represents the image embedding function, which maps an image to a point in Euclidean space, as outlined in [5]. While peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM) serve as common image metrics, it is necessary to note that these metrics primarily operate at the per-pixel level, failing to capture differences in semantics and human perception.

Deep learning-based image similarity metrics have the capacity to capture semantics to a certain extent. Johnson et al. [33] introduced two concepts known as *perceptual loss functions*, which enable the measurement of high-level perceptual and semantic distinctions between images. These perceptual loss functions are constructed using a loss network denoted as  $\phi$ , which is pre-trained for image classification. It is worth noting that these perceptual loss functions themselves are deep convolutional neural networks. The perceptual loss functions consist of a feature reconstruction loss and a style reconstruction loss. Let  $\phi_j(A)$  be the activations of the  $j$ -th layer of the network  $\phi$  when processing the image  $A$ , and let  $L$  represent the shape. Then the *feature reconstruction loss* is the Euclidean distance between feature representations:

$$\mathcal{L}_{\text{feat}}(A, B) = \frac{1}{L} \|\phi(A) - \phi(B)\|_2^2. \quad (12)$$

The *style reconstruction loss* is responsible for capturing the stylistic characteristic of images. It is defined as the squared Frobenius norm of the difference between the Gram matrices,  $G_l^\phi$ , of two images, and it is expressed as follows:

$$\mathcal{L}_{\text{style}}(A, B) = \|G_l^\phi(A) - G_l^\phi(B)\|_F^2. \quad (13)$$

Deep features have proven to be highly effective in semantic tasks and serve as robust models for understanding human perceptual behavior. Notably, Zhang et al. [34] conducted a comprehensive evaluation of deep features across various architectural designs and tasks. Their research compared these deep features with traditional metrics, and the results demonstrated a significant superiority of deep features. They outperformed previous metrics by substantial margins, particularly on a newly introduced dataset focused on human perceptual similarity judgments.

In a related development, Wang et al. [35] proposed a deep ranking model designed to learn fine-grained image similarity models. It utilizes a triplet-based hinge loss ranking function to characterize fine-grained image similarity relationships. It also incorporates a multiscale neural network architecture capable of capturing both global visual properties and image semantics. Additionally, in 2023, Zhu et al. [36] proposed ViTScore, a novel semantic similarity evaluation metric for images. ViTScore relies on the pre-trained image model ViT (Vision Transformer) and represents a cutting-edge approach to assessing semantic similarity in the context of images.

2) **Text:** In the context of text transmission, conventional metrics such as word-error rate (WER) often struggle to effectively address semantic tasks, as pointed out by Farsad et al. [37]. In response to this challenge, the bilingual evaluation understudy (BLEU) metric, initially designed for machine translation evaluations by Papineni et al. [38], has found utility in the domain of semantic communication. Specifically, BLEU assesses the quality of semantic communication as follows: Let  $l_a$  and  $l_b$  represent the word lengths of sentences  $a$  and  $b$ , respectively, then the BLEU score is defined by

$$\log \text{BLEU} = \min(1 - \frac{l_a}{l_b}, 0) + \sum_{n=1}^N w_n \log p_n, \quad (14)$$

where  $w_n$  is the weights of the  $n$ -grams, and  $p_n$  denotes the  $n$ -grams score, which is defined as

$$p_n = \frac{\sum_k \min(C_k(\hat{s}, C_k(s)))}{\sum_k \min(C_k(\hat{s}))}, \quad (15)$$

where  $C_k(\cdot)$  is the frequency count function for the  $k$ -th element in the  $n$ -th gram.

The concept of sentence similarity, as proposed in [39], serves as a metric for quantifying the semantic similarity between two sentences. It is expressed as:

$$\tau(\hat{s}, s) = \frac{\mathbf{B}_\Phi(s) \cdot \mathbf{B}_\Phi(\hat{s})^T}{\|\mathbf{B}_\Phi(s)\| \|\mathbf{B}_\Phi(\hat{s})\|}. \quad (16)$$

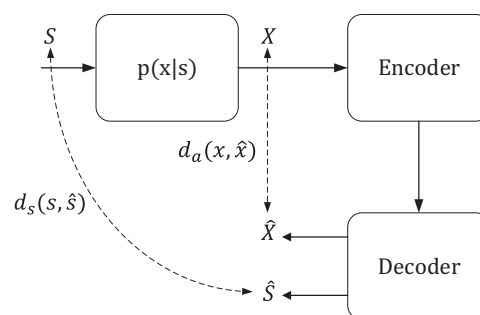
where  $\mathbf{B}_\Phi(\cdot)$  represents the BERT model [40], which maps a sentence to its semantic vector space. This model is pre-trained on a massive dataset comprising billions of sentences, enabling it to capture rich semantic information.

3) **Audio:** In the realm of semantic communication, novel perception-based audio metrics are exploited, including perceptual evaluation of speech quality (PESQ) [41], short-time objective intelligibility (STOI) [42], perceptual objective listening quality assessment (POLQA) [43] and unconditional kernel deep speech distance (KDS) [44]. These metrics provide valuable insights into the semantic aspects of audio quality and perception.

#### 4.2. Semantic Rate-Distortion Theorem

In the context of semantic communication, the process of feature extraction and coding representation at the semantic level plays a crucial role in reducing information redundancy and extracting the most salient semantic features, thus improving the effectiveness of semantic transmission. The semantic rate-distortion theorem is a theoretical framework designed to address the challenges associated with distortion and encoding in semantic communication. It offers solutions and insights into optimizing the trade-off between preserving semantic content and achieving efficient encoding. Liu et al. [45] have introduced a comprehensive semantic rate distortion theory framework. In this framework, they consider a memoryless information source represented as a tuple of random variables, denoted as  $(S, X)$ . It has a joint probability distribution denoted as  $p(s, x)$  within a product alphabet  $\mathcal{S} \times \mathcal{X}$ . Here,  $S$  represents the intrinsic state, capturing the "semantic" aspect of the source, which is not directly observable. On the other hand,  $X$  represents the extrinsic observation of the source, capturing the "appearance" as perceived by an observer.

For a length- $n$  independent and identically distributed (i.i.d.) sequence from the source, denoted as  $(S^n, X^n)$ , a source encoder  $f_n$  with a rate of  $R$  is a mapping that transforms  $X^n$  into an index within the set  $\{1, 2, \dots, 2^{nR}\}$ . This encoder corresponds to a decoder that maps the index back into a pair of sequences, denoted as  $(\hat{S}^n, \hat{X}^n)$ , where these sequences are drawn from product alphabet  $\hat{\mathcal{S}} \times \hat{\mathcal{X}}$ . This process is illustrated in Figure 3.



**Figure 3.** A semantic communication system for introducing semantic rate distortion[45].

In this framework, two distortion metrics are considered:  $d_s(s, \hat{s})$  representing semantic distortion, and  $d_a(x, \hat{x})$  representing appearance distortion. These metrics map elements from alphabets  $\mathcal{S} \times \hat{\mathcal{S}}$  and  $\mathcal{X} \times \hat{\mathcal{X}}$  to non-negative real numbers. Consequently, the block-wise distortions are defined as:

$$d_s(s^n, \hat{s}^n) = \frac{1}{n} \sum_{i=1}^n d_s(s_i, \hat{s}_i) \quad (17)$$

$$d_a(x^n, \hat{x}^n) = \frac{1}{n} \sum_{i=1}^n d_a(x_i, \hat{x}_i) \quad (18)$$

Moreover, the framework defines the semantic rate distortion function as

$$R(D_S, D_a) = \min I(X; \hat{S}, \hat{X}), \quad (19)$$

$$\text{s.t. } \mathbb{E} \hat{d}_s(X, \hat{S}) \leq D_s, \quad \mathbb{E} \hat{d}_a(X, \hat{X}) \leq D_a \quad (20)$$

where  $S$  and  $\hat{S}$  represent the semantic understanding of the sender and the receiver, while  $X$  and  $\hat{X}$  are their respective semantic representations. Expanding on this, Guo et al. [46] delved into the analysis of semantic rate-distortion involving two users, considering the perspective of side information. This perspective can be expressed as:

$$R(D_1, D_2, D_s) = \min I(X_1, X_2; \hat{X}_1, \hat{X}_2, \hat{S} | Y), \quad (21)$$

where  $X_1$  and  $X_2$  are the semantic representations of two users, respectively.  $Y$  represents the side information. In 2022, Stavrou and Kountouris [47] further studied the characteristics of this system, particularly focusing on the Hamming distortion metric.

### 4.3. Semantic Coding

Semantic rate-distortion theory directly corresponds to coding technology. For a given transmitting task, a semantic coding strategy needs to achieve two potentially conflicting goals:

- Maximizing expected faithfulness (minimizing expected semantic distortion).
- Minimizing expected coding length.

An ideal semantic coding strategy should simultaneously minimize both expected semantic distortion and expected coding length. However, achieving this delicate balance is highly complex and challenging. In current practices, a common approach involves the use of a dual distortion metric to represent semantic coding. Shao et al. [48] used semantic distortion and semantic cost to define the achievable region for semantic coding. Semantic distortion reflects the disparities in semantic understanding between the receiver and the sender. Semantic cost, which quantifies the simplicity or understandability of information, is often represented as the length of the corresponding bit sequence. The definition of the achievable distortion-cost regions can be expressed as: A distortion-cost pair  $(L, D)$  is achievable if there exists a semantic encoding scheme  $U$  if  $D_U = D, L_U = L$ .

Agheli [49] et al. have explored semantic coding within a multi-user context. They introduced an approach where observations from an information source are filtered and sent to two monitors depending on their importance for each user's specific objectives. By optimizing codeword lengths using semantics-aware utility functions, substantial reductions in the amount of communicated status updates can be achieved. Xiao et al. [50] proposed the rate-distortion theory of strategic semantic communication. Their approach integrates game theory models with rate-distortion theory to characterize how information interaction between semantic encoders and decoders impacts communication distortion performance. Furthermore, Tang et al. [51] considered a semantic source that consists of a set of correlated random variables whose joint probabilistic distribution can be

described by a Bayesian network. Their work focuses on characterizing the limits of lossy compression for semantic sources and establishing upper and lower bounds for the rate-distortion function.

## 5. Semantic Channel Coding

Channel capacity is the most successful and central contribution to Shannon's information theory. On the one hand, it provides the maximum number of distinguishable signals through repeated use of a communication channel. By appropriately mapping the source into "widely spaced" input sequences for the channel, one can transmit a message with an exceedingly low probability of error, subsequently reconstructing the source message at the output. On the other hand, channel capacity represents the rate at which reliable information can be transmitted through a noisy channel, as discussed by Verdú in 'Fifty Years of Shannon Theory' [52].

Similarly, the issue of capacity holds immense significance within the realm of semantic communication. In this section, we delve into the concepts of semantic noise and semantic channel capacity. Furthermore, several widely concerned questions about semantic channel capacity are raised and addressed. Finally, we attempt to give a general description of semantic channel capacity.

In the domain of digital communications, a discrete channel is defined as a system comprising an input alphabet denoted by  $\mathcal{X}$ , an output alphabet represented by  $\mathcal{Y}$ , and a probability transition matrix  $p(y|x)$  that quantifies the probability of observing the output symbol  $y$  when transmitting the message  $x$ .

**Definition 2.** *The channel capacity of a discrete memoryless channel is defined as*

$$C = \max_{p(x)} I(X; Y), \quad (22)$$

where the maximum is taken over all possible input distributions  $p(x)$  provided channel transition probability function  $p(y|x)$ .

For the sake of convenient comparison, we commonly refer to this as the physical channel capacity. Moreover, Shannon's theorem established that information can be transmitted reliably through a physical channel at any rate up to the channel capacity, known as the channel coding theorem.

**Theorem 2.** (Channel coding theorem) *For a discrete memoryless channel, all rates below capacity  $C$  are achievable. Specifically, for every rate  $R < C$ , there exists a sequence of  $(2^{nR}, n)$  codes with maximum probability of error  $\lambda^{(n)} \rightarrow 0$ , as  $n \rightarrow \infty$ . Conversely, any sequence of  $(2^{nR}, n)$  codes with  $\lambda^{(n)} \rightarrow 0$ , as  $n \rightarrow \infty$ , must have  $R \leq C$ .*

The channel coding theorem states that all rates below capacity  $C$  are achievable, while rates exceeding this capacity are unachievable. This leads us to contemplate the significance and formulation of channel capacity in the context of semantic communications, a topic of interest among scholars in the field. In the exploration of semantic information theory, we try to address the following three fundamental questions regarding capacity.

1. Is there an analogous concept of channel capacity in semantic communications, which we may term 'semantic channel capacity'?
2. Is it possible that the semantic channel capacity is greater than the physical channel capacity?
3. Is there a universal expression for semantic channel capacity?

Next, we will address these three fundamental questions and introduce the semantic noise along with the semantic channel capacity theorem.



5.1. Semantic Noise

Noise introduces uncertainty into communication systems and also poses challenges to communication technologies. In the absence of noise, the transmission and exchange of any information are perfect and lossless, rendering capacity a meaningless concept. Generally, semantic noise exists widely in semantic communications. Prior to the formal introduction, it is essential to clarify that semantic noise and semantic channel noise are distinct concepts. In most scholarly literature, semantic noise refers to the mismatch of semantics. Semantic channel noise commonly refers to the discrepancies in the knowledge background of both parties in semantic communications. It is worth noting that the semantic noise may be added either at the physical channel or at the semantic channel.

In their respective works, Qin et al. [5], Shi et al. [7], and Hu et al. [53] offered similar definitions of semantic noise within the context of semantic communication. Qin et al. defined semantic noise as a disturbance that affects the interpretation of a message, characterized by a semantic information mismatch between the transmitter and receiver. Similarly, Shi et al. described semantic noise as noise introduced during the communication process, leading to misunderstanding and incorrect reception of semantic information. It can be introduced at various stages, including encoding, data transportation, and decoding. Hu et al. defined semantic noise as a unique form of noise in semantic communication systems, denoting the misalignment between intended semantic symbols and received ones.

Semantic noise varies across different source categories. The semantic noise in text refers to semantic ambiguity, which slightly changes the semantic meaning of a sentence. In the case of images, it can be modeled by using adversarial samples. The following examples illustrate instances of semantic ambiguity in the text, categorized by communication channel [15]:

- The meaning of a message is changed due to transmission errors, e.g., from "contend" to "content" (Physical channel)
- Translation of one natural language into another language where some concepts in the two languages have no precise match (Semantic channel)
- Communicating parties use different background knowledge to understand the message (e.g., Apple has different meanings in vegetable markets and mobile phone stores) (Semantic channel)

An illustrative example of an adversarial image is presented in Figure 4, in which the adversarial samples are added. It is apparent that the image when perturbed with adversarial noise, can mislead deep learning models in their classification, while remaining visually indistinguishable from the original image to human observers [54].

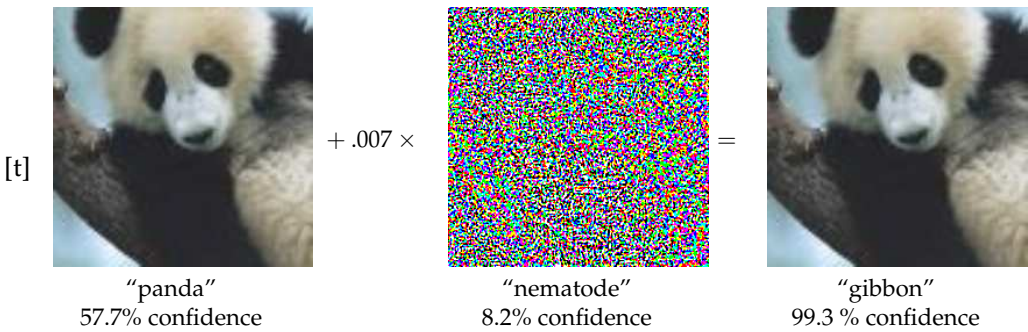


Figure 4. An example of the adversarial sample in the image [54].

In summary, semantic noise encompasses both semantic channel noise and physical channel noise. It represents a discrepancy in semantic information between the transmitter and the receiver within the context of semantic communication.

5.2. Semantic Capacity

For the first question in Section 5, we argue that determining the capacity of the semantic channel is a challenging task, and expressing it using current mathematical tools remains elusive. On the one

hand, to the best of our knowledge, it is still an open problem to effectively model the local knowledge and global knowledge shared by both communicating parties. On the other hand, quantifying and representing semantic channel noise is also difficult. However, there are some studies that delve into semantic capacity. We would like to clarify that the semantic channel capacity studied in most current works is the capacity at the semantic level, rather than the capacity of the semantic channel. In the subsequent sections of this paper, the semantic channel capacity we refer to is also the information capacity at the semantic level.

In 2016, Okamoto [55] argued that the semantic channel capacity represents the maximum rate of semantic information that can be transmitted over the semantic channel, or the ratio of the maximum semantic communication amount to the communication data size. Similarly, Bao et al. [15] defined the semantic channel capacity as the capacity limit such that a transmission rate can be achieved with arbitrarily small semantic errors within the limit. Specifically, they derived the semantic channel coding theorem.

**Theorem 3.** (Semantic Channel Coding Theorem I) For every discrete memoryless channel, the channel capacity

$$C_s = \sup_{P(X|Z)} \{I(X;Y) - H(Z|X) + \overline{H_S(Y)}\} \quad (23)$$

has the following property: For any  $\epsilon > 0$  and  $R \leq C_s$ , there is a block coding strategy such that the maximal probability of semantic error is not greater than  $\epsilon$ .

Among them,  $X$  and  $Y$  serve as the input and output of the channel, while  $Z$  is the semantic representation.  $I(X;Y)$  denotes the mutual information between  $X$  and  $Y$ .  $H(Z|X)$  represents the semantic uncertainty associated with the encoding. Additionally,  $\overline{H_S(Y)}$  represents the average logical information of the received message.

Based on the theorem presented above, we can see that semantic capacity may be higher or lower than the physical channel capacity, depending on whether  $\overline{H_S(Y)}$  or  $H(Z|X)$  is larger. This observation implies that through the utilization of a semantic encoder with minimal semantic ambiguity and a semantic decoder possessing robust inference capabilities or an extensive knowledge base, it is possible to achieve high-rate semantic communication using a low-rate physical channel.

It is worth noting that in semantic communication, even when some symbolic errors occur, the underlying semantics may still remain unchanged. In other words, semantic communications may allow a certain level of bit errors, signified by a non-zero BER. Consequently, in cases where the communication system permits a given non-zero BER ( $P_{e,b} \geq 0$ ), the transmission rate  $R$  can exceed the physical channel capacity  $C$ . As a response to the second question, Figure 5 shows the error probability of the communication system. It indicates that when  $R > C$ , the  $P_{e,b}$  becomes greater than zero, while  $P_{e,s}$  may still remain at zero.

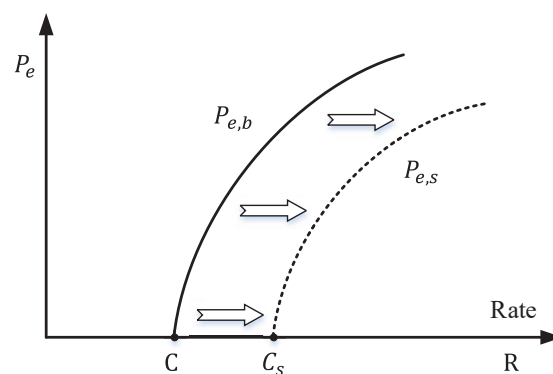


Figure 5. Error probability of the communication system.

Consider a semantic communication system, as illustrated in Figure 6. In this system, a semantic sender aims to reliably transmit a latent semantic message  $S$  within the message  $W$  at a rate  $R$  bits per transmission to a semantic receiver over a noisy physical channel. The source message set is defined as  $[1 : 2^{nR}] = \{1, 2, \dots, 2^{nR}\}$ , which contains a semantic message subset  $[1 : 2^{\lceil \alpha nR \rceil}] = \{1, 2, \dots, 2^{\lceil \alpha nR \rceil}\}$ , where the coefficient  $\alpha$  falls within the range  $0 \leq \alpha < 1$ . Thus, given the semantic mapping  $[1 : 2^{nR}] \rightarrow [1 : 2^{\alpha nR}]$  and the discrete memoryless channel  $p(y|x)$ , Ma et al. [56] define the semantic channel  $C_s$  as follows:

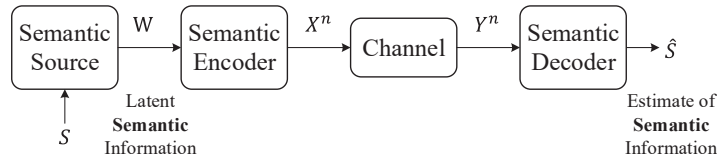


Figure 6. A semantic communication system

$$C_s = \max_{p(x)} \frac{I(X;Y)}{\alpha}. \quad (24)$$

Moreover, they further give the following theorem concerning semantic channel coding.

**Theorem 4.** (Semantic Channel Coding Theorem II) For every bit rate  $R < C_s = \max_{p(x)} \frac{I(X;Y)}{\alpha}$ , there exists a sequence of  $(2^{nR}, n)$  codes with average probability of error  $P_{e,s}^{(n)}$  that tends to zero as  $n \rightarrow \infty$ .

Conversely, for every sequence of  $(2^{nR}, n)$  codes with probability of error  $P_{e,s}^{(n)}$  that tends to zero as  $n \rightarrow \infty$ , the rate must satisfy  $R \leq C_s = \max_{p(x)} \frac{I(X;Y)}{\alpha}$ .

For the third question, we believe that the semantic channel capacity should be related to the specific task and the background knowledge possessed by both parties, in other words, it is task and goal-oriented. Additionally, we contend that this semantic channel capacity should be dynamic, adapting to the temporal relevance of information. Consequently, establishing a universally applicable expression for semantic channel capacity becomes a complex undertaking. However, if we try to describe it, we propose that it possesses three distinct characteristics.

1. It takes the form of conditions. These conditions at least shall include tasks, public knowledge, private knowledge, and goals.
2. It has the capability to reflect the temporal value of information. For instance, in situations demanding low latency, messages with slow transmission rates will possess a low information value density.
3. It should encompass the concept of physical channel capacity since the semantic channel does not really exist, and the transmission of symbols must still be achieved through the real physical channel.

## 6. Related Mathematical Theories and Methods

In this section, we will introduce some mathematical tools and concepts that are highly relevant to semantic communication. These include the Age of Information, Information Bottleneck, Large Language Models, and Joint Source-Channel Coding. In fact, these methods have already been extensively applied within the realm of semantic communication and have demonstrated their effectiveness. They are expected to become an auxiliary or part of semantic information theory.

### 6.1. Age of Information (AoI)

In the context of semantic communication, the transmitted semantics are usually task-oriented. Many of these tasks require low latency and are sensitive to the freshness of the message. Regardless

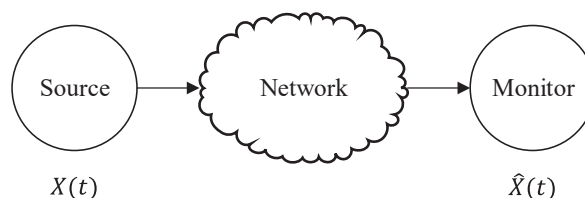
of how accurately the message can be recovered, if it arrives too late, it can be rendered completely useless, as highlighted by Gündüz in their work [2]. For example, in the realm of autonomous driving, vehicles depend on near-instantaneous data exchange to make split-second decisions, thereby avoiding potential traffic accidents. Similarly, in the context of remote surgery systems, the timely update of surgical tool positions is necessary to ensure the safety and precision of medical procedures. However, this does not imply that the transmitter should update the current state as rapidly as possible. When the network is congested, such a strategy could lead to a monitor receiving delayed updates that were backlogged in the communication system.

Conversely, it is important to note that neither the source statistics nor quality measures change over time in digital communications. Therefore, there exists a necessity to acquire the status of remote sensors or systems. In other words, semantic communication needs to take into account the temporal aspect and its impact on the overall effectiveness of the communication process.

The Age of Information (AoI) serves as an end-to-end performance metric, providing a means to quantify the timeliness of a monitor's knowledge about a particular entity or process [57]. It has the potential to empower semantic communication to measure the freshness of semantics. We will discuss several kinds of AoI and the relationship that exists between semantic communication and AoI.

As depicted in Figure 7, a source continuously generates new updates to a network that are subsequently delivered to a destination monitor. In this context, the source is represented as a random process denoted as  $X(t)$ , while the monitor possesses the capability to estimate the current state  $\hat{X}(t)$ . Each update packet is associated with a timestamp, denoted as  $u$ , and its age at time  $t \geq u$  is defined as  $t - u$ . An update is said to be fresh when its timestamp matches the current time  $t$ , resulting in an age of zero [58]. When the monitor's most recently received update at time  $t$  has a timestamp of  $u(t)$ , the age of information is defined as

$$\Delta(t) = t - u(t). \quad (25)$$



**Figure 7.** Updates from a source pass through the network to a destination monitor.

Figure 8 depicts a visualization of AoI at the monitor over time. In this scenario, the transmitter sends update packets according to a First-Come-First-Served (FCFS) queuing discipline, allowing only one packet transmission at any given time. Since the monitor sees updates that are delivered at times  $t'_j$  after traveling through the network, its age process  $\Delta t$  is reset to  $\Delta(t'_j) = t'_j - t_j$ , which is the age of update  $j$  when it is delivered. Building upon this foundation, various representations of age can be derived, such as time average age and peak age.

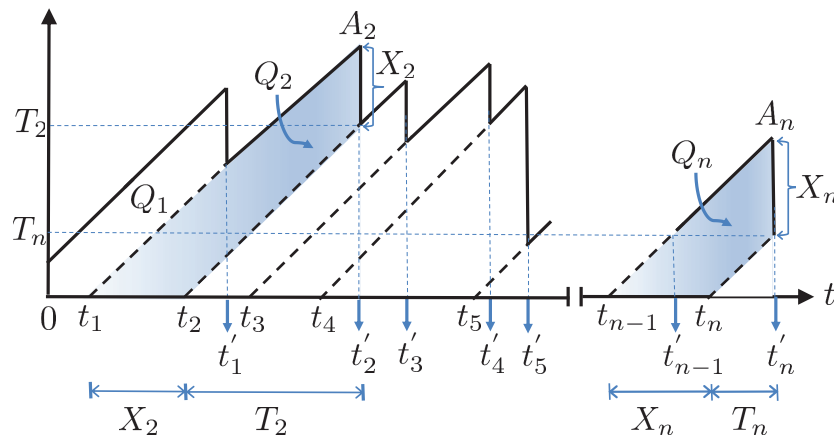


Figure 8. AoI evolution vs. time for  $n$  update packets [59].

For an interval of observation  $(0, T)$ , the time average age of a status update system is defined as

$$\Delta_T = \frac{1}{T} \int_0^T \Delta(t) dt. \quad (26)$$

The time average age can be represented as the area under  $\Delta(t)$ , which corresponds to the shaded region denoted as  $Q_n$  in Figure 8. It is calculated as

$$Q_n = \frac{1}{2} (T_n + X_n)^2 - \frac{1}{2} T_n^2 = X_n T_n + X_n^2 / 2. \quad (27)$$

When  $(X_n, T_n)$  is a stationary ergodic process, the time average age satisfies

$$\Delta_T = \frac{\mathbb{E}[Q_n]}{\mathbb{E}[T_n]} = \frac{\mathbb{E}[X_n T_n] + \mathbb{E}[X_n^2] / 2}{\mathbb{E}[T_n]}. \quad (28)$$

On the other hand, the difficulty in evaluating  $\mathbb{E}[X_n T_n]$  prompted the introduction of peak age of information (PAoI) [60], an alternative and more manageable metric for age assessment. In this context, let  $X_i$  be the interarrival time of the  $i$ th update, and  $T_i$  be the corresponding system time. Then, the *Peak Age of Information* (PAoI) metric is defined as the value of age achieved immediately before receiving the  $i$ th update

$$A_i = X_i + T_i. \quad (29)$$

Moreover, the average peak age of a status update system is calculated as

$$A_T = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^{N(T)} A_i, \quad (30)$$

where  $N(T)$  is the number of samples that completed by time  $T$ . The peak age can be utilized in applications where the worst-case age is of interest or where there is a need to apply a threshold restriction on age.

While AoI proves valuable in assessing the freshness of information and has seen widespread application across various system contexts employing diverse queuing disciplines, it falls short in effectively capturing the informational content of transmitted packets and the monitor's current knowledge. In fact, even when the monitor has perfect knowledge of the underlying process in question, AoI continuously increases over time, leading to unnecessary penalties. This motivated the introduction of a novel metric, known as the age of incorrect information (AoII), which addresses the limitations of conventional AoI and error penalty functions within the framework of status updates [61]. AoII is defined as



$$\Delta(t) = (t - v(t)) \cdot \mathbb{1}\{X(t) \neq \hat{X}(t)\}, \quad (31)$$

where  $\mathbb{1}$  is the indicator function. When  $\mathbb{1}\{X(t) \neq \hat{X}(t)\} = 0$ , the monitor has the most updated information about  $X(t)$  irrespective of when the status update was received. AoII extends the notion of fresh updates and adequately captures the information content that the updates bring to the monitor.

In a given scenario, a transmitter observes the source process  $X(t)$  and sends samples/updates regarding this source over time to a monitor. The primary objective is to reconstruct the original process,  $X(t)$ , at the monitor utilizing the received samples/updates [62]. In this general setup, the transmitter has two key decisions to take during each time slot: i) whether to sample  $X(t)$  or not and ii) whether to transmit the available or newly generated sample/update or not. The real-time compression/transmission has been explored in the literature, with relevant studies such as [63,64].

Moreover, the objective function that takes into account semantic awareness can be formulated as:

$$\mathcal{S} = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T c(t), \quad (32)$$

where  $c(t)$  represents a cost function that is selected appropriately. For instance,  $c(t)$  can be chosen as

$$c(t) = a(t) \cdot \mathbb{1}\{X(t) \neq \hat{X}(t)\}, \quad (33)$$

where  $a(t)$  is a function of  $t$ . When one does not take channel or transmission delays into account, the problem described above can be viewed as a rate-distortion problem. Nevertheless, the primary objective in this context is to devise sampling and transmission policies that minimize the average distortion with consideration for delay sensitivity.

In their work, Gündüz et al. [2] represented three distinct sampling and transmission policies that have subtle connections between the age metrics and semantics of information: AoI-aware, source-aware, semantics-aware sampling and transmission policy. The third policy addresses the limitations of the prior ones by considering not only the state of the source signal but also the status of the reconstructed signals at the monitor. It is expected to be introduced into the semantic communication systems to reflect the impact of time on semantics.

In semantic communication, combining AoI/AoII and semantic distortion can be a way to evaluate delay-sensitive distortion. In [65], it constructs a multi-user uplink non-orthogonal multiple access system to analyze its transmission performance by harnessing the age of incorrect information. It adopts the semantic similarity [39] as the performance metric and AoII as the time delay, thus, the instantaneous AoII of the  $k$ -th user in the scenario can be expressed as

$$\Psi^k(t) \triangleq \Delta_{\text{AoII}}(X_t, \hat{X}_t, t) \quad (34)$$

$$= (1 - \psi_n^k(s, \hat{s})) \cdot \Delta(t) \quad (35)$$

$$= (1 - \frac{\mathbf{B}_\Phi(s) \cdot \mathbf{B}_\Phi(\hat{s})^T}{\|\mathbf{B}_\Phi(s)\| \|\mathbf{B}_\Phi(\hat{s})\|})(t - u(t)). \quad (36)$$

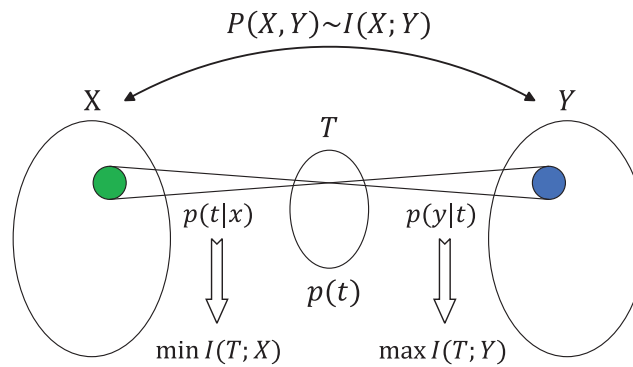
By minimizing the average cost of the system, the optimization problem of multiple users can be solved. This system utilizes AoII as a metric while simultaneously assessing semantic similarity and AoI performance.

## 6.2. Information Bottleneck (IB)

Since Deep Neural Networks (DNNs) have demonstrated excellent performance in extracting features and key information, many semantic communication systems use them as a tool adopted in encoders and decoders. Information bottleneck (IB), an information-theoretic principle, can explore

the interplay between the informativeness of extracted features and its effect on communication goals. It offers a novel paradigm for analyzing Deep Neural Networks (DNNs) and tries to shed light on their layered structure, generalization capabilities, and learning dynamics [66], which has drawn attention from machine learning and information theory [67]. The interaction of IB and DNNs in the literature can be divided into two main categories. The first is to use the IB theories in order to analyze DNNs and the other is to use the ideas from IB to improve the DNN-based learning algorithms [68]. Therefore, the information bottleneck theory may also be a powerful mathematical tool for the development of semantic communication.

The IB method was introduced as an information-theoretic principle for extracting relevant information from an input random variable  $X \in \mathcal{X}$  with respect to an output random variable  $Y \in \mathcal{Y}$  [69]. In the information bottleneck framework, the objective is to compress the information that variable  $X$  carries about variable  $Y$  through a compact 'bottleneck' represented as  $T$ . Either way, the basic trade-off is between minimizing the compression information and maximizing the relevant information [70]. An illustration of this idea is given in Figure 9.



**Figure 9.** The information between  $X$  and  $Y$  is squeezed through the compact "bottleneck" representation,  $T$ . In particular, under some constraint over the minimal level of relevant information,  $I(T; Y)$ , one is trying to minimize the compression information,  $I(T; X)$ .

Given a joint probability distribution  $p(x, y)$  the IB optimization problem can be formulated as follows: find  $T$  such that the mutual information  $I(T; X)$  is minimized, subject to the constraint that  $I(T; Y)$  does not exceed a predetermined threshold  $\hat{D}$ . Consequently, it is intuitive to introduce a mathematical function analogous to the rate-distortion function

$$\hat{R}(\hat{D}) \triangleq \min_{\{p(t|x): I(T; Y) \leq \hat{D}\}} I(T; X). \quad (37)$$

In a word,  $\hat{R}(\hat{D})$  represents the minimal achievable compression information, for which the relevant information is above  $\hat{R}(\hat{D})$ . Moreover, the optimal assignment can be determined by minimizing the corresponding functional

$$\mathcal{L} = I(T; X) - \beta I(T; Y), \quad (38)$$

where  $\beta$  is the Lagrange multiplier attached to the constrained meaningful information while maintaining the normalization of the mapping  $p(t; x)$  for every  $x$ . At  $\beta = 0$  the quantization is the most sketchy possible, everything is assigned to a single point, while as  $\beta \rightarrow \infty$  we are pushed toward arbitrarily detailed quantization. By varying the parameter  $\beta$ , one can explore the tradeoff between preserving meaningful information and achieving compression at different levels of resolutions [71].

When applying the IB method to enhance semantic communication systems, a common approach involves utilizing IB to extract task-related information while eliminating redundant features. Barbarossa et al. [72] presented an approach to semantic communication building on the information bottleneck principle. It used the information bottleneck as a way to identify relevant information and

adapt the information bottleneck online, as a function of the wireless channel state, in order to strike an optimal trade-off between transmit power, reconstruction accuracy, and delay. More specifically, Barbarossa et al. considered the case where the transmitted data are corrupted by noise, so that the receiver does not have direct access to  $T$ , but only to a corrupted version  $T + \eta$ , where  $\eta$  is the channel noise. Then, they redefined the bottleneck optimization problem as

$$\min_{A, M} I(X; T) + \beta \cdot \text{MSE}(Y, \hat{Y}), \quad (39)$$

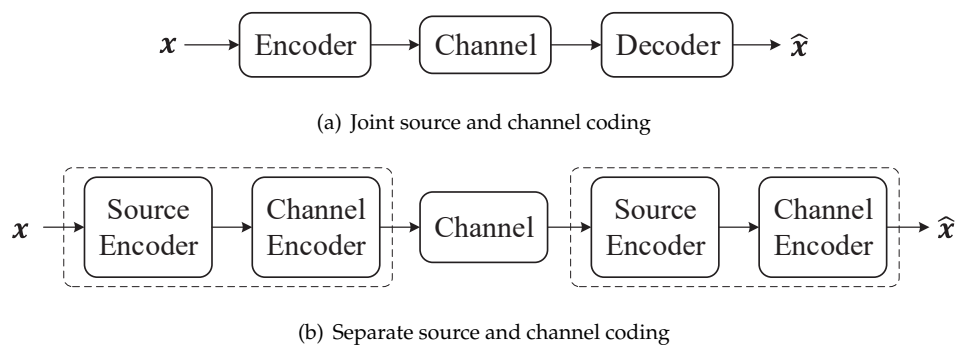
where  $T = A \cdot X + \xi$  denotes a linear encoder, including additive noise  $\xi$ , while  $\hat{Y} = M \cdot (T + \eta)$  is a linear estimate of  $Y$ ;  $\beta$  is a scalar parameter that allows it to tune the trade-off between complexity and relevant information: small values of  $\beta$  lead to small complexity encoders, but possibly large distortion. Conversely, larger values of  $\beta$  lead to reduced distortion at the expense of increased complexity.

Li et al. [73] used the information bottleneck framework to extract distinguishable features in-distribution data while keeping their compactness and informativeness. Wei et al. [74] presented a federated semantic learning framework to collaboratively train the semantic-channel encoders of multiple devices with the coordination of a base station-based semantic-channel decoder. In this approach, the information bottleneck is leveraged to drive the loss design by formalizing a rate-distortion tradeoff. This tradeoff serves to eliminate the redundancies of semantic features while maintaining task-relevant information.

In 2018, Zaslavsky et al. [75] presented empirical evidence that IB may give rise to human-like semantic representations. They conducted research into how human languages categorize colors. Their findings indicated that human languages appear to evolve under pressure to efficiently compress meanings into communication signals by optimizing the information bottleneck tradeoff between informativeness and complexity. Furthermore, Tucker et al. [76] studied how trading off three factors -utility, informativeness, and complexity - shapes emergent communication, including compared to human communication. Their study not only shed light on these factors' impact on communication but also made comparisons to human communication processes.

### 6.3. Joint Source Channel Coding

From the perspective of structural design in communication systems, two primary categories emerge: separate source and channel coding, and joint source and channel coding. In the case of joint source and channel coding, as depicted in Figure 10(a), there exists only an encoder and decoder. In this setup, the system optimizes both source coding and channel coding. These coding schemes are integrated into a unified process. Conversely, Figure 10(b) illustrates the separate source and channel coding system. These two design approaches differ significantly. In Shannon's theory, source codes are designed independently to achieve efficient data representation, while channel codes are designed separately, tailored to the specific characteristics of the channel.



**Figure 10.** The joint/separate source and channel coding system.

Shannon's separation theorem is a fundamental result in information theory. It established that separate encoders can achieve the same rates as the joint encoder. Specifically, it tied together the two basic theorems of information theory: data compression and data transmission, as outlined in [32]. The data compression theorem is an outcome of the Asymptotic Equipartition Property (AEP), which shows that there exists a "small" subset (of size  $2^{nH}$ ) of all possible source sequences that contain most of the probability. Consequently, one can represent the source with a small probability of error using an average of  $H$  bits per symbol. The data transmission theorem, on the other hand, is based on the joint AEP. It capitalizes on the fact that for long block lengths, the output sequence of the channel is very likely to be jointly typical with the input codeword, while any other codeword is jointly typical with probability  $\sim 2^{-nI}$ . As a result, we can employ approximately  $2^{nI}$  codewords while still having a negligible probability of error. The source-channel separation theorem shows that one can design the source code and the channel code separately and combine the results to achieve the same optimal performance.

Shannon's separation theorem proves that the two-step source and channel coding approach is theoretically optimal in the asymptotic limit of infinitely long source and channel blocks. However, in practical applications, it is widely recognized that joint source and channel coding surpasses the separate approach in terms of performance especially for limited source and channel coding complexity. Furthermore, JSCC is resilient to variations in channel conditions and does not suffer from abrupt quality degradations, commonly referred to as the "cliff effect" in digital communication systems [77]. Nonetheless, semantic communication is often oriented towards emerging intelligent applications from the Internet-of-Things to autonomous driving and tactile Internet requires transmission of image/video data under extreme latency, bandwidth, and energy constraints. This precludes computationally demanding long-blocklength source and channel coding techniques. Therefore, joint source and channel coding may be a potential trend in semantic communication systems.

In recent years, semantic communication systems, employing joint source and channel coding, have demonstrated superior performance across various domains, surpassing separate design approaches and offering new potential applications and advantages [78,79]. In 2019, a noteworthy development known as Deep JSCC [77] introduced encoder and decoder functions parameterized by two convolutional neural networks, trained jointly. The results show that the deep JSCC scheme outperforms digital transmission concatenating JPEG or JPEG2000 compression with a capacity-achieving channel code at low SNR and channel bandwidth values in the presence of additive white Gaussian noise (AWGN). Building upon this foundation, Deep JSCC-f [80] investigated how noiseless or noisy channel output feedback can be incorporated into the Deep JSCC to improve the reconstruction quality at the receiver. Xie et al. [39] introduced DeepSC, a deep learning-based semantic communication system designed for text transmission. In comparison to traditional communication systems that do not account for semantic information exchange, DeepSC exhibits remarkable robustness to channel variations and superior performance, particularly in low SNR conditions. For speech signals, DeepSC-S [81] was designed, which outperforms traditional communications in both cases in terms of the speech signals metrics, such as signal-to-distortion ratio and perceptual evaluation of speech distortion. Lastly, MU-DeepSC [82] represents a multi-user semantic communication system for transmitting multimodal data. Its transceiver is ingeniously designed and optimized to capture features from correlated multimodal data, facilitating task-oriented transmission. These recent advancements highlight the growing potential and versatility of joint source and channel coding in semantic communication systems across a range of data types and applications.

#### 6.4. Large Language Models

In a semantic communication system, the knowledge base plays a crucial role in enabling the semantic encoder and decoder to comprehend and infer semantic information. In general, knowledge includes public knowledge and private knowledge. The former is shared by all communication participants, while the latter is unique to a user. In fact, the knowledge base is a key feature that

distinguishes semantic communication from conventional communication systems. However, the representation and updating of knowledge is a challenging task, which is also one of the factors that make semantic communication difficult to mathematically model. In recent years, Large Language Models (LLMs) developed rapidly and have shown great potential in intelligent tasks [83]. In this subsection, we will introduce the feasibility of introducing LLMs into semantic communication. We believe that LLMs may play a potential role in knowledge bases for semantic communication.

Language models (LMs) are computational models that have the capacity to understand and generate human language [84]. These models possess the transformative ability to predict the likelihood of word sequences or generate new text based on a given input [85]. In the context of a sequence denoted as  $X$ , LM tasks aim to predict the next token, denoted as  $y$ . The model is trained by maximizing the probability of the given token sequence conditioned on the context, i.e.,  $P(y|X) = P(y|x_1, x_2, \dots, x_{n-1})$ , where  $x_1, x_2, \dots, x_{n-1}$  are the tokens in the context sequence, and  $n$  is the current position. Utilizing the chain rule, the conditional probability can be decomposed into a product of probabilities at each position:

$$p(y|X) = \prod_{n=1}^N P(y_n|x_1, x_2, \dots, x_{n-1}), \quad (40)$$

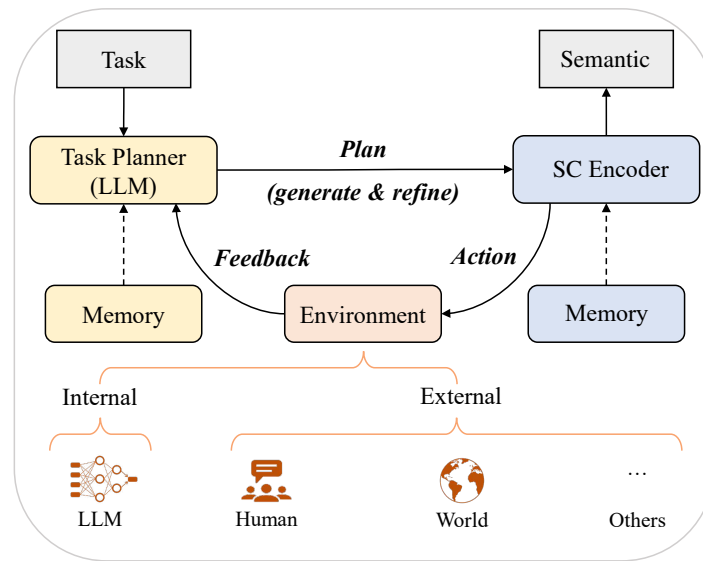
where  $N$  represents the length of the sequence. Consequently, the model predicts each token at each position in an autoregressive manner, ultimately generating a complete text sequence.

Large Language Models (LLMs) are advanced language models with massive parameter sizes and exceptional learning capabilities. The core module behind many LLMs such as GPT-3 [86], InstructGPT [87], and GPT-4 is the self-attention module in the Transformer [88] architecture. This self-attention module serves as the foundational building block for various language modeling tasks.

A fundamental characteristic of Large Language Models (LLMs) lies in their capability for in-context learning. This means the model is trained to generate text based on a provided context or prompt. This capability empowers LLMs to produce responses that are not only more coherent but also contextually relevant, making them well-suited for interactive and conversational applications. Another crucial aspect of LLMs is Reinforcement Learning from Human Feedback (RLHF) [89]. This technique involves fine-tuning the model by using human-generated responses as rewards, enabling the model to learn from its mistakes and progressively enhance its performance over time.

Since LLMs use a large amount of data, parameters, and even human feedback during training, they have a perception and understanding of the world, which can be called the knowledge base to a certain extent. On the other hand, due to the excellent performance and capabilities of LLMs, they have the potential to be applied in a variety of intelligent tasks. Zhao et al. [90] formulated the general planning paradigm of LLMs for solving complex tasks. We think it provides insights into semantic communication based on a large language model, which is illustrated in Figure 11.





**Figure 11.** An illustration of the formulation for prompt-based planning by LLMs for semantic communication (transmitter)[90].

Figure 11 illustrates the process by which LLMs help the transmitter's semantic encoder extract semantics. In this paradigm, there are typically three components: task planner, semantic encoder, and environment. Specifically, task planner, which is played by LLMs, aims to generate the whole plan to solve a target task-oriented communication. The plan can be presented in various forms, e.g. a visual question answering task [82] or a text transmission task [39]. Then, the semantic encoder is responsible for executing the actions in the plan and generating semantics. It can be implemented by models like LLMs for textual tasks. Furthermore, environment refers to where the semantic encoder generates the semantics, which can be set differently according to specific tasks. It provides feedback about the execution result of the action to the task planner, either in the form of natural language or from other multimodal signals.

In current semantic communication systems, the construction of the knowledge base (KB) faces several issues, including limited knowledge representation, frequent knowledge updates, and insecure knowledge sharing. Jiang et al. [91] proposed using LLMs as knowledge bases in semantic communication. Specifically, they represented three LLMs-based KBs in SC models, which is illustrated in Figure 12. 1) GPT-based KBs. In the context of text-based SC systems, the KB must possess the ability to comprehend textual content and identify various subjects, their attributes, and relationships. ChatGPT, an AI assistant developed by OpenAI based on the GPT-3.5 model, demonstrates an accurate understanding of text content and can provide correct responses to a wide array of questions. Utilizing ChatGPT as the KB for textual data allows for the extraction of key content from input text, tailored to user requirements. 2) SAM-based KBs. For image-based SC systems, the KB should be capable of segmenting various objects within an image and recognizing their respective categories and relationships. An intriguing AI model suited for this purpose is the Segment Anything Model (SAM), introduced by Meta AI [92]. SAM is a groundbreaking segmentation system with the unique ability to generalize zero-shot to unfamiliar images and objects without requiring additional training. As such, SAM can effectively serve as the KB for image-related tasks. 3) WavLM-Based KBs. To enable SC systems for audio, the KB must encompass a wide range of audio-related tasks, including automatic speech recognition, speaker identification, and speech separation. WavLM [93], a large-scale audio model proposed by Microsoft Research Asia (MSRA), stands as a potential solution for such applications.

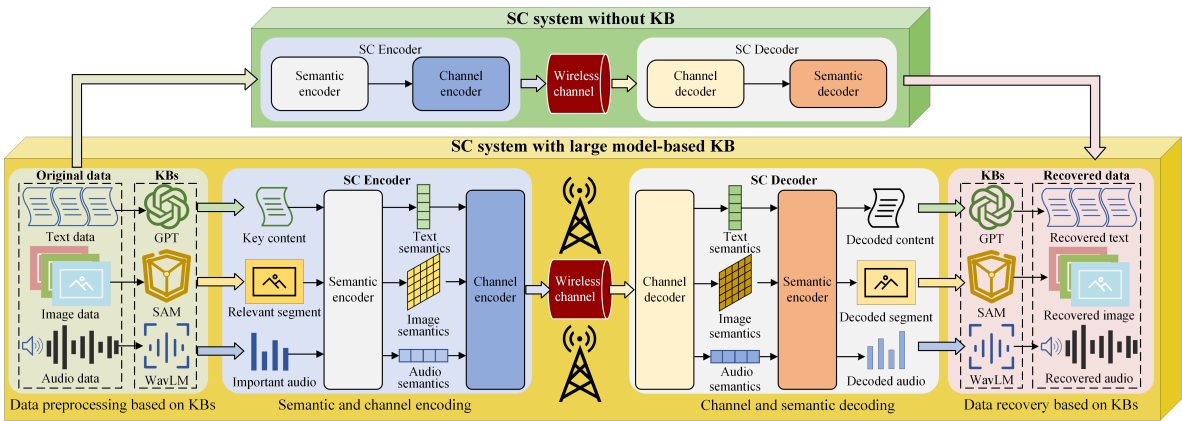


Figure 12. Implementation of large AI models-based KBs in different SC models [91].

7. Challenges

In this section, we present several challenges that arise in the field of semantic information theory and semantic communication. The mathematical theory of semantic communication and the mathematical representation of semantics are referred to as the semantic information theory. Distinguishing semantic information theory from classical information theory reveals several notable disparities:

- Whether a message is true or not is irrelevant in classical information theory.
- Whether a message is related to the task/time is indifferent in classical information theory.
- Whether a message can effectively convey meaning is not a concern of classical information theory.

However, these differences or concerns about semantic communication are challenging issues at a theoretical level. In fact, the development of semantic information theory is in the initial stage, with a large number of open problems that have not yet been solved, such as.

1. **The Role of Semantics in Data Compression and Reliable Communication:** How can semantics contribute to data compression and enhance the reliability of communication?
2. **Relationship Between Semantic and Engineering Coding:** What is the interplay between semantic coding/decoding techniques and conventional engineering coding/coding problems?
3. **Fundamental Limits of Semantic Communication:** Are there established limits or boundaries in semantic coding?
4. **Enhancing Efficiency and Reliability in Semantic Communication:** What factors should be taken into account to improve efficiency and reliability in semantic communication?
5. **Principles for DL-Based Semantic Communication:** How should we architect the framework of a semantic communication system rooted in deep learning, and what theoretical guidance exists?
6. **Capacity of Semantic-Aware Networks:** What is the capacity of a semantic network, and how can we evaluate the performance limits of a semantic transmission network?

Beyond the theoretical realm, semantic communication systems present an array of open challenges:

1. **Scheduling and Energy Optimization:** Delving into scheduling and resource allocation policies within semantic communications, with a concentrated effort on optimizing energy utilization.
2. **Complexity of Semantic-Enabled Networks:** Semantic-enabled networks face high complexity due to the need to share knowledge with users. This necessitates a framework for evaluating the complexity and necessity of semantic communication networks.

3. **Multi-criteria optimization:** Developing strategies for semantic communication in scenarios where multiple tasks and objectives coexist.
4. **Knowledge Updates Tracking:** Recognizing that knowledge can evolve over time within semantic networks.
5. **Applications:** Identifying specific use cases and applications that align with semantic communication systems.
6. **Performance Metrics:** Defining comprehensive performance metrics for assessing the effectiveness and efficiency of semantic communication systems.

## 8. Concluding Remarks

Semantic communication, as an innovative communication structure, has revolutionized the traditional data transmission paradigm and holds the potential to provide fresh insights into large-scale intelligent processing services. Nevertheless, it is important to note that the field of semantic communication is still in its infancy, offering abundant opportunities for further exploration and research.

This article systematically summarized the development of semantic information theory, encompassing a comprehensive review of related advancements. Beginning with an exploration of semantic entropy, we further delved into statistical probability, logical probability, semantic rate distortion, semantic encoding, semantic noise, and semantic channel capacity. Moreover, the article presented a structured exposition of the mathematical theories and methodologies relevant to semantic communication, including concepts like the age of information, information bottleneck, and joint source-channel coding.

In addition, we investigated the prevalent challenges and open problems within the realm of semantic information theory and semantic communication. We believe that this article will make a meaningful contribution to the complete establishment of semantic information theory and the rapid evolution of semantic communication.

**Author Contributions:** Conceptualization, P.F. and G.X.; methodology, P.F. and K.B.L.; software, G.X.; writing—original draft preparation, G.X.; writing—review and editing, P.F. and K.B.L.; visualization, P.F. and K.B.L.; supervision, P.F. and K.B.L.; All authors have read and agreed to the published version of the manuscript.

**Funding:** This work was supported by the National Key Research and Development Program of China (Grant NO. 2021.YFA1000504)

**Institutional Review Board Statement:** Not applicable

**Informed Consent Statement:** Not applicable

**Data Availability Statement:** Not applicable

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. Letaief, K.B.; Chen, W.; Shi, Y.; Zhang, J.; Zhang, Y.J.A. The roadmap to 6G: AI empowered wireless networks. *IEEE communications magazine* **2019**, *57*, 84–90.
2. Gündüz, D.; Qin, Z.; Aguerri, I.E.; Dhillon, H.S.; Yang, Z.; Yener, A.; Wong, K.K.; Chae, C.B. Beyond transmitting bits: Context, semantics, and task-oriented communications. *IEEE Journal on Selected Areas in Communications* **2022**, *41*, 5–41.
3. Yang, W.; Du, H.; Liew, Z.Q.; Lim, W.Y.B.; Xiong, Z.; Niyato, D.; Chi, X.; Shen, X.S.; Miao, C. Semantic communications for future internet: Fundamentals, applications, and challenges. *IEEE Communications Surveys & Tutorials* **2022**.
4. Weaver, W. Recent contributions to the mathematical theory of communication. *ETC: a review of general semantics* **1953**, pp. 261–281.
5. Qin, Z.; Tao, X.; Lu, J.; Tong, W.; Li, G.Y. Semantic communications: Principles and challenges. *arXiv preprint arXiv:2201.01389* **2021**.

6. Strinati, E.C.; Barbarossa, S. 6G networks: Beyond Shannon towards semantic and goal-oriented communications. *Computer Networks* **2021**, *190*, 107930.
7. Shi, G.; Xiao, Y.; Li, Y.; Xie, X. From semantic communication to semantic-aware networking: Model, architecture, and open problems. *IEEE Communications Magazine* **2021**, *59*, 44–50.
8. Kountouris, M.; Pappas, N. Semantics-empowered communication for networked intelligent systems. *IEEE Communications Magazine* **2021**, *59*, 96–102.
9. Kalfa, M.; Gok, M.; Atalik, A.; Tegin, B.; Duman, T.M.; Arikan, O. Towards goal-oriented semantic signal processing: Applications and future challenges. *Digital Signal Processing* **2021**, *119*, 103134.
10. Lan, Q.; Wen, D.; Zhang, Z.; Zeng, Q.; Chen, X.; Popovski, P.; Huang, K. What is semantic communication? A view on conveying meaning in the era of machine intelligence. *Journal of Communications and Information Networks* **2021**, *6*, 336–371.
11. Uysal, E.; Kaya, O.; Ephremides, A.; Gross, J.; Codreanu, M.; Popovski, P.; Assaad, M.; Liva, G.; Munari, A.; Soret, B.; et al. Semantic communications in networked systems: A data significance perspective. *IEEE Network* **2022**, *36*, 233–240.
12. Zhang, P.; Xu, W.; Gao, H.; Niu, K.; Xu, X.; Qin, X.; Yuan, C.; Qin, Z.; Zhao, H.; Wei, J.; et al. Toward wisdom-evolutionary and primitive-concise 6G: A new paradigm of semantic communication networks. *Engineering* **2022**, *8*, 60–73.
13. Shi, Y.; Zhou, Y.; Wen, D.; Wu, Y.; Jiang, C.; Letaief, K.B. Task-oriented communications for 6g: Vision, principles, and technologies. *arXiv preprint arXiv:2303.10920* **2023**.
14. Shannon, C.E. A mathematical theory of communication. *The Bell system technical journal* **1948**, *27*, 379–423.
15. Bao, J.; Basu, P.; Dean, M.; Partridge, C.; Swami, A.; Leland, W.; Hendler, J.A. Towards a theory of semantic communication. In Proceedings of the 2011 IEEE Network Science Workshop. IEEE, 2011, pp. 110–117.
16. Iyer, S.; Khanai, R.; Torse, D.; Pandya, R.J.; Rabie, K.M.; Pai, K.; Khan, W.U.; Fadlullah, Z. A survey on semantic communications for intelligent wireless networks. *Wireless Personal Communications* **2023**, *129*, 569–611.
17. Carnap, R. *Logical foundations of probability*; Vol. 2, Citeseer, 1962.
18. Carnap, R.; Bar-Hillel, Y.; et al. An outline of a theory of semantic information **1952**.
19. Floridi, L. Outline of a theory of strongly semantic information. *Minds and machines* **2004**, *14*, 197–221.
20. D'Alfonso, S. On quantifying semantic information. *Information* **2011**, *2*, 61–101.
21. Basu, P.; Bao, J.; Dean, M.; Hendler, J. Preserving quality of information by using semantic relationships. *Pervasive and Mobile Computing* **2014**, *11*, 188–202.
22. Chattopadhyay, A.; Haeffele, B.D.; Geman, D.; Vidal, R. Quantifying task complexity through generalized information measures **2020**.
23. Melamed, I.D. Measuring semantic entropy. In Proceedings of the Tagging Text with Lexical Semantics: Why, What, and How?, 1997.
24. Liu, X.; Jia, W.; Liu, W.; Pedrycz, W. AFSSE: An interpretable classifier with axiomatic fuzzy set and semantic entropy. *IEEE Transactions on Fuzzy Systems* **2019**, *28*, 2825–2840.
25. De Luca, A.; Termini, S. A definition of a nonprobabilistic entropy in the setting of fuzzy sets theory. In *Readings in Fuzzy Sets for Intelligent Systems*; Elsevier, 1993; pp. 197–202.
26. Choi, J.; Loke, S.W.; Park, J. A unified view on semantic information and communication: A probabilistic logic approach. In Proceedings of the 2022 IEEE International Conference on Communications Workshops (ICC Workshops). IEEE, 2022, pp. 705–710.
27. Xin, G.; Fan, P. EXK-SC: A semantic communication model based on information framework expansion and knowledge collision. *Entropy* **2022**, *24*, 1842.
28. Kolchinsky, A.; Wolpert, D.H. Semantic information, autonomous agency and non-equilibrium statistical physics. *Interface focus* **2018**, *8*, 20180041.
29. Rényi, A. On measures of entropy and information. In Proceedings of the Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics. University of California Press, 1961, Vol. 4, pp. 547–562.
30. Venhuizen, N.J.; Crocker, M.W.; Brouwer, H. Semantic entropy in language comprehension. *Entropy* **2019**, *21*, 1159.
31. Lu, C. From Bayesian inference to logical Bayesian inference: A new mathematical frame for semantic communication and machine learning. In Proceedings of the Intelligence Science II: Third IFIP TC 12

- International Conference, ICIS 2018, Beijing, China, November 2-5, 2018, Proceedings 2. Springer, 2018, pp. 11–23.
32. Cover, T.M.; Thomas, J.A. *Elements of information theory*; 1991.
  33. Johnson, J.; Alahi, A.; Fei-Fei, L. Perceptual losses for real-time style transfer and super-resolution. In Proceedings of the Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14. Springer, 2016, pp. 694–711.
  34. Zhang, R.; Isola, P.; Efros, A.A.; Shechtman, E.; Wang, O. The unreasonable effectiveness of deep features as a perceptual metric. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 586–595.
  35. Wang, J.; Song, Y.; Leung, T.; Rosenberg, C.; Wang, J.; Philbin, J.; Chen, B.; Wu, Y. Learning fine-grained image similarity with deep ranking. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2014, pp. 1386–1393.
  36. Zhu, T.; Peng, B.; Liang, J.; Han, T.; Wan, H.; Fu, J.; Chen, J. How to Evaluate Semantic Communications for Images with ViTScore Metric? *arXiv preprint arXiv: 2309.04891* **2023**.
  37. Farsad, N.; Rao, M.; Goldsmith, A. Deep learning for joint source-channel coding of text. In Proceedings of the 2018 IEEE international conference on acoustics, speech and signal processing (ICASSP). IEEE, 2018, pp. 2326–2330.
  38. Papineni, K.; Roukos, S.; Ward, T.; Zhu, W.J. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the Proceedings of the 40th annual meeting of the Association for Computational Linguistics, 2002, pp. 311–318.
  39. Xie, H.; Qin, Z.; Li, G.Y.; Juang, B.H. Deep learning enabled semantic communication systems. *IEEE Transactions on Signal Processing* **2021**, *69*, 2663–2675.
  40. Devlin, J.; Chang, M.W.; Lee, K.; Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* **2018**.
  41. Rix, A.W.; Beerends, J.G.; Hollier, M.P.; Hekstra, A.P. Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs. In Proceedings of the 2001 IEEE international conference on acoustics, speech, and signal processing. Proceedings (Cat. No. 01CH37221). IEEE, 2001, Vol. 2, pp. 749–752.
  42. Taal, C.H.; Hendriks, R.C.; Heusdens, R.; Jensen, J. An algorithm for intelligibility prediction of time–frequency weighted noisy speech. *IEEE Transactions on Audio, Speech, and Language Processing* **2011**, *19*, 2125–2136.
  43. Beerends, J.G.; Schmidmer, C.; Berger, J.; Obermann, M.; Ullmann, R.; Pomy, J.; Keyhl, M. Perceptual objective listening quality assessment (polqa), the third generation itu-t standard for end-to-end speech quality measurement part i—temporal alignment. *journal of the audio engineering society* **2013**, *61*, 366–384.
  44. Bińkowski, M.; Donahue, J.; Dieleman, S.; Clark, A.; Elsen, E.; Casagrande, N.; Cobo, L.C.; Simonyan, K. High fidelity speech synthesis with adversarial networks. *arXiv preprint arXiv:1909.11646* **2019**.
  45. Liu, J.; Zhang, W.; Poor, H.V. A rate-distortion framework for characterizing semantic information. In Proceedings of the 2021 IEEE International Symposium on Information Theory (ISIT). IEEE, 2021, pp. 2894–2899.
  46. Guo, T.; Wang, Y.; Han, J.; Wu, H.; Bai, B.; Han, W. Semantic compression with side information: A rate-distortion perspective. *arXiv preprint arXiv:2208.06094* **2022**.
  47. Stavrou, P.A.; Kountouris, M. A rate distortion approach to goal-oriented communication. In Proceedings of the 2022 IEEE International Symposium on Information Theory (ISIT). IEEE, 2022, pp. 590–595.
  48. Shao, Y.; Cao, Q.; Gunduz, D. A theory of semantic communication. *arXiv preprint arXiv:2212.01485* **2022**.
  49. Agheli, P.; Pappas, N.; Kountouris, M. Semantic Source Coding for Two Users with Heterogeneous Goals. In Proceedings of the GLOBECOM 2022-2022 IEEE Global Communications Conference. IEEE, 2022, pp. 4983–4988.
  50. Xiao, Y.; Zhang, X.; Li, Y.; Shi, G.; Başar, T. Rate-distortion theory for strategic semantic communication. In Proceedings of the 2022 IEEE Information Theory Workshop (ITW). IEEE, 2022, pp. 279–284.
  51. Tang, J.; Yang, Q.; Zhang, Z. Information-Theoretic Limits on Compression of Semantic Information. *arXiv preprint arXiv:2306.02305* **2023**.
  52. Verdu, S. Fifty years of Shannon theory. *IEEE Transactions on information theory* **1998**, *44*, 2057–2078.



53. Hu, Q.; Zhang, G.; Qin, Z.; Cai, Y.; Yu, G.; Li, G.Y. Robust semantic communications with masked VQ-VAE enabled codebook. *IEEE Transactions on Wireless Communications* **2023**.
54. Goodfellow, I.J.; Shlens, J.; Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* **2014**.
55. Okamoto, T. A unified paradigm of organized complexity and semantic information theory. *arXiv preprint arXiv:1608.00941* **2016**.
56. Ma, S.; Wu, Y.; Qi, H.; Li, H.; Shi, G.; Liang, Y.; Al-Dhahir, N. A Theory for Semantic Communications. *arXiv preprint arXiv:2303.05181* **2023**.
57. Kosta, A.; Pappas, N.; Angelakis, V.; et al. Age of information: A new concept, metric, and tool. *Foundations and Trends® in Networking* **2017**, *12*, 162–259.
58. Yates, R.D.; Sun, Y.; Brown, D.R.; Kaul, S.K.; Modiano, E.; Ulukus, S. Age of information: An introduction and survey. *IEEE Journal on Selected Areas in Communications* **2021**, *39*, 1183–1210.
59. Abd-Elmagid, M.A.; Pappas, N.; Dhillon, H.S. On the role of age of information in the Internet of Things. *IEEE Communications Magazine* **2019**, *57*, 72–77.
60. Costa, M.; Codreanu, M.; Ephremides, A. Age of information with packet management. In Proceedings of the 2014 IEEE International Symposium on Information Theory. IEEE, 2014, pp. 1583–1587.
61. Maatouk, A.; Kriouile, S.; Assaad, M.; Ephremides, A. The age of incorrect information: A new performance metric for status updates. *IEEE/ACM Transactions on Networking* **2020**, *28*, 2215–2228.
62. Sun, Y.; Polyanskiy, Y.; Uysal, E. Sampling of the Wiener process for remote estimation over a channel with random delay. *IEEE Transactions on Information Theory* **2019**, *66*, 1118–1135.
63. Witsenhausen, H.S. On the Structure of Real-Time Source Coders. *Bell System Technical Journal* **1979**, *58*, 1437–1451.
64. Mahajan, A.; Teneketzis, D. Optimal design of sequential real-time communication systems. *IEEE Transactions on Information Theory* **2009**, *55*, 5317–5338.
65. Chen, J.; Wang, J.; Jiang, C.; Wang, J. Age of Incorrect Information in Semantic Communications for NOMA Aided XR Applications. *IEEE Journal of Selected Topics in Signal Processing* **2023**.
66. Goldfeld, Z.; Polyanskiy, Y. The information bottleneck problem and its applications in machine learning. *IEEE Journal on Selected Areas in Information Theory* **2020**, *1*, 19–38.
67. Shao, J.; Mao, Y.; Zhang, J. Learning task-oriented communication for edge inference: An information bottleneck approach. *IEEE Journal on Selected Areas in Communications* **2021**, *40*, 197–211.
68. Hafez-Kolahi, H.; Kasaei, S. Information bottleneck and its applications in deep learning. *arXiv preprint arXiv:1904.03743* **2019**.
69. Tishby, N.; Zaslavsky, N. Deep learning and the information bottleneck principle. In Proceedings of the 2015 IEEE information theory workshop (itw). IEEE, 2015, pp. 1–5.
70. Slonim, N. The information bottleneck: Theory and applications. PhD thesis, Hebrew University of Jerusalem Jerusalem, Israel, 2002.
71. Tishby, N.; Pereira, F.C.; Bialek, W. The information bottleneck method. *arXiv preprint physics/0004057* **2000**.
72. Barbarossa, S.; Comminiello, D.; Grassucci, E.; Pezone, F.; Sardellitti, S.; Di Lorenzo, P. Semantic communications based on adaptive generative models and information bottleneck. *arXiv preprint arXiv:2309.02387* **2023**.
73. Li, H.; Yu, W.; He, H.; Shao, J.; Song, S.; Zhang, J.; Letaief, K.B. Task-Oriented Communication with Out-of-Distribution Detection: An Information Bottleneck Framework. *arXiv preprint arXiv:2305.12423* **2023**.
74. Wei, H.; Ni, W.; Xu, W.; Wang, F.; Niyato, D.; Zhang, P. Federated Semantic Learning Driven by Information Bottleneck for Task-Oriented Communications. *IEEE Communications Letters* **2023**.
75. Zaslavsky, N.; Kemp, C.; Regier, T.; Tishby, N. Efficient human-like semantic representations via the information bottleneck principle. *arXiv preprint arXiv:1808.03353* **2018**.
76. Tucker, M.; Shah, J.; Levy, R.; Zaslavsky, N. Towards human-agent communication via the information bottleneck principle. *arXiv preprint arXiv:2207.00088* **2022**.
77. Bourtsoulatz, E.; Kurka, D.B.; Gündüz, D. Deep joint source-channel coding for wireless image transmission. *IEEE Transactions on Cognitive Communications and Networking* **2019**, *5*, 567–579.
78. Zhang, H.; Shao, S.; Tao, M.; Bi, X.; Letaief, K.B. Deep learning-enabled semantic communication systems with task-unaware transmitter and dynamic data. *IEEE Journal on Selected Areas in Communications* **2022**, *41*, 170–185.



79. Xie, H.; Qin, Z.; Tao, X.; Letaief, K.B. Task-oriented multi-user semantic communications. *IEEE Journal on Selected Areas in Communications* **2022**, *40*, 2584–2597.
80. Kurka, D.B.; Gündüz, D. DeepJSCC-f: Deep joint source-channel coding of images with feedback. *IEEE Journal on Selected Areas in Information Theory* **2020**, *1*, 178–193.
81. Weng, Z.; Qin, Z. Semantic communication systems for speech transmission. *IEEE Journal on Selected Areas in Communications* **2021**, *39*, 2434–2444.
82. Xie, H.; Qin, Z.; Li, G.Y. Task-oriented multi-user semantic communications for VQA. *IEEE Wireless Communications Letters* **2021**, *11*, 553–557.
83. Huang, J.; Chang, K.C.C. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403* **2022**.
84. Min, B.; Ross, H.; Sulem, E.; Veyseh, A.P.B.; Nguyen, T.H.; Sainz, O.; Agirre, E.; Heintz, I.; Roth, D. Recent advances in natural language processing via large pre-trained language models: A survey. *ACM Computing Surveys* **2023**, *56*, 1–40.
85. Chang, Y.; Wang, X.; Wang, J.; Wu, Y.; Zhu, K.; Chen, H.; Yang, L.; Yi, X.; Wang, C.; Wang, Y.; et al. A survey on evaluation of large language models. *arXiv preprint arXiv:2307.03109* **2023**.
86. Floridi, L.; Chiriatti, M. GPT-3: Its nature, scope, limits, and consequences. *Minds and Machines* **2020**, *30*, 681–694.
87. Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* **2022**, *35*, 27730–27744.
88. Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, Ł.; Polosukhin, I. Attention is all you need. *Advances in neural information processing systems* **2017**, *30*.
89. Christiano, P.F.; Leike, J.; Brown, T.; Martic, M.; Legg, S.; Amodei, D. Deep reinforcement learning from human preferences. *Advances in neural information processing systems* **2017**, *30*.
90. Zhao, W.X.; Zhou, K.; Li, J.; Tang, T.; Wang, X.; Hou, Y.; Min, Y.; Zhang, B.; Zhang, J.; Dong, Z.; et al. A survey of large language models. *arXiv preprint arXiv:2303.18223* **2023**.
91. Jiang, F.; Peng, Y.; Dong, L.; Wang, K.; Yang, K.; Pan, C.; You, X. Large AI Model-Based Semantic Communications. *arXiv preprint arXiv:2307.03492* **2023**.
92. Kirillov, A.; Mintun, E.; Ravi, N.; Mao, H.; Rolland, C.; Gustafson, L.; Xiao, T.; Whitehead, S.; Berg, A.C.; Lo, W.Y.; et al. Segment anything. *arXiv preprint arXiv:2304.02643* **2023**.
93. Chen, S.; Wang, C.; Chen, Z.; Wu, Y.; Liu, S.; Chen, Z.; Li, J.; Kanda, N.; Yoshioka, T.; Xiao, X.; et al. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing* **2022**, *16*, 1505–1518.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.