

# Less Data, More Knowledge: Building Next Generation Semantic Communication Networks

Christina Chaccour, *Graduate Student Member, IEEE*, Walid Saad, *Fellow, IEEE*,  
Mérrouane Debbah, *Fellow, IEEE*, Zhu Han, *Fellow, IEEE*, and H. Vincent Poor, *Life Fellow, IEEE*

**Abstract**—Semantic communication is viewed as a revolutionary paradigm that can potentially transform how we design and operate wireless communication systems. However, despite a recent surge of research activities in this area, remarkably, the research landscape is still limited in at least three ways. First, the very definition of a “semantic communication system” remains ambiguous, and it differs from one work to another. Second, there is a lack of *fundamental and scalable* frameworks for building next-generation semantic communication networks based on rigorous and well-defined technical foundations. Third, the question of what a “semantic representation” means, and on how this representation can be used to instill meaning, significance, and structure to every information transfer over a wireless network remain unanswered. In this tutorial, we present the first rigorous and holistic vision of an end-to-end semantic communication network that is founded on novel concepts from artificial intelligence (AI), causal reasoning, transfer learning, and minimum description length theory. We first discuss how the design of semantic communication networks requires a move from *data-driven* and *information-driven* AI-augmented networks, in which wireless networks remain “tied” to data, towards *knowledge-driven* and *reasoning-driven* AI-native networks in which wireless networks are AI-native and can perform versatile logic. We then distinguish the concept of semantic communications from several other approaches that have been conflated with it. For instance, we opine that effectively and efficiently building next-generation semantic communication networks must go beyond: a) creating a new type of encoder and decoder at the transmitter/receiver side, and b) designing a new “AI for wireless” framework in which AI is used to extract some application features or to fine tune a wireless protocol or algorithm. Then, we identify the main tenets that are needed to build an end-to-end semantic communication network. Among those building blocks of a semantic communication networks, we highlight the necessity of creating semantic representations of data that satisfy the key properties of *minimalism, generalizability, and efficiency* so as to faithfully represent the data and enable the transmitter and receiver to do more with less, i.e., computationally generate content via a minimally semantic representation. We then explain how those representations can form the basis a so-called *semantic language* that will allow a transmitter and receiver to communicate at a semantic level. In this regard, we distinguish the concept of a semantic language from that of a natural language, and we

present the pillars needed to gradually build a semantic language with fundamental structural content, yet tolerable complexity. We then show that, by using semantic representation and languages, the traditional transmitter and receiver now become a teacher and apprentice. The teacher can identify the semantic content elements in the raw datastream and learn its semantic representation. The apprentice can reason over a semantic representation, map its corresponding semantic content element, and further draw logical conclusions based on the cumulative knowledge base built. This phenomenon mimics the growth of a child’s language’s expressivity and reasoning in a more-or-less parallel fashion. We then concretely define the concept of *reasoning* by investigating the fundamentals of causal representation learning and their role in designing reasoning-driven semantic communication networks. We particularly demonstrate that reasoning faculties are majorly characterized by the ability to capture causal and associational relationships in datastreams. This enables radio nodes to communicate minimal, generalizable, and efficient semantic representations, and ultimately perform versatile logical conclusions – doing more with less. For such reasoning-driven networks, we revisit the fundamentals of information theory, in order to emphasize the concepts that must be redefined to capture semantic reasoning. We then propose novel and essential semantic communication *key performance indicators (KPIs)* and metrics that include new “reasoning capacity” measures that could go beyond Shannon’s bound to capture the imminent convergence of computing and communication resources. Finally, we explain how semantic communications can be scaled to large-scale networks such as cellular networks (6G and beyond), and deployed in emerging environments such as open radio access networks (O-RAN). In a nutshell, we expect this tutorial to provide a unified and self-contained reference on how to properly build, design, analyze, and deploy next-generation semantic communication networks.

**Index Terms**— Semantic communications, Semantic language, Causality, Knowledge, Reasoning, 6G, AI-Native, Machine Learning, Beyond 6G.

## I. INTRODUCTION

Future wireless systems, namely 6G systems and beyond, must cater to the complex and stringent requirements of emerging applications such as the metaverse, holographic teleportation, digital twins, and Industry 5.0 [1]. Nonetheless, delivering a disruptive leap in wireless technologies cannot be fulfilled by continuing to pursue incremental advances to conventional wireless system components such as spectrum and multi-antenna technologies. Instead, it is necessary to rethink the way in which the entire wireless system architecture and functions are designed and operated. Along those line, current 5G and 6G research efforts, have already demonstrated the efficiency of using AI-driven augmentation

C. Chaccour and W. Saad are with Wireless@VT, Bradley Department of Electrical and Computer Engineering, Virginia Tech, Arlington, VA, USA, Emails: christinac@vt.edu, walids@vt.edu.

M. Debbah is with the Technology Innovation Institute, 9639 Masdar City, Abu Dhabi, United Arab Emirates and also with CentraleSupélec, University Paris-Saclay, 91192 Gif-sur-Yvette, France, Email: merouane.debbah@tii.ae.

Z. Han is with the Department of Electrical and Computer Engineering, University of Houston, Houston, TX 77004 USA, and also with the Department of Computer Science and Engineering, Kyung Hee University, Seoul 446-701, South Korea.

H.V. Poor is with Department of Electrical and Computer Engineering, Princeton University, NJ 08544, USA, poor@princeton.edu.

in addressing various challenges throughout the network layer stack. For example, in [2]–[4], it was reported that artificial intelligence (AI) and machine learning (ML) can provide robust, accurate, and reduced complexity for tasks like channel estimation, initialization, and symbol detection. Also, ML has played an integral role routing protocol design [5], resource management [6], and network management [7], among many other specific wireless problems and protocols [8].

However, current AI-based wireless approaches [2]–[8] remain limited in a number of ways. First, in the state-of-art, the design of wireless networks and protocols is either limited to *data-driven solutions* or to *information-driven approaches*, and, thus, it fails to *leverage knowledge* accumulated throughout the operation of the system, as shown in Fig. 1. For instance, under the data-driven paradigm, wireless networks rely on discrete elements, (e.g., spectrum data, channel data, quality-of-service (QoS) values) to fine-tune their operation. Here, *wireless networks designs, and their performance will be tied very closely to data*. Meanwhile, in information-based approaches, wireless networks leverage *information-centric* metrics, such as age of information (AoI), value-of-information, application data, and reliability to make more informed network decisions (which links multiple data points into an important performance-evaluating metric). Nonetheless, both of these approaches fail to leverage, accumulate, and organize *knowledge within the data and information sets*. Essentially, the decision making performed via these two paradigms remains largely training dependent, and exhibits *limited generalizability*. In contrast, under our *envisioned knowledge-driven* and *reasoning-driven* (see Fig. 1) paradigms, the wireless network will be able to proactively make decisions and draw conclusions on its own based on the logical intent that can be extracted from built knowledge bases. In this case, the network will be able to acquire *more knowledge with less data*, compared to data/information-driven approaches. This, in turn, will facilitate the network’s ability to achieve high-rate, low latency, and high-reliability, thereby becoming more adept in meeting the very stringent QoS requirements of future 6G and beyond applications. In addition, by continuously exploiting accumulated knowledge, the network can now reduce reliance on the brittle spectral resources, enhance the way in which data is transmitted and recovered, and rely on computational resources to yield semantic content rather than reconstruct raw data. We envision that these two paradigms must be a stepping stone of the emerging concept of *AI-native* wireless systems, that has attracted significant attention by academia, industry, and standardization bodies.

Broadly, the concept of AI-native systems envisions building the entire protocol stack and air-interface of a wireless system using AI techniques. For instance, in [9] a new research direction is attempting to transforming the air-interface into a full AI-AI integration, to reach the milestone of AI-nativeness. However, existing works [9]–[11] in this regard, do not specify what type of AI framework should be used to build such

systems, but they still hint towards classical techniques (e.g., convolutional neural networks, reinforcement learning, etc). In this regard, we opine that AI-native networks cannot continue to rely on mere AI-augmented techniques, such as designing a transceiver via autoencoders or resource management protocols via deep reinforcement learning (RL) [2]–[8]. In contrast, instead of *augmenting* or merely *replacing* existing network layers and components with AI to achieve AI-nativeness, we envision that, ultimately, AI-native wireless systems should be intrinsically designed and structured based on *applied knowledge*. In our envisioned *reasoning-driven AI native wireless systems*, the end-to-end (E2E) design and operation of the network will be designed using next-generation reasoning AI frameworks that can exploit causality and stochasticity in the data, identify structure, deduce logical connections, and extract (and mitigate) semantic noise (the source, root-cause, and mechanism of the noise are identified). In this scenario, the wireless system becomes a *living, sustainable network* that can grow with its cumulative knowledge in order to execute operations that cannot simply be done when being overly reliant on existing data and re-training mechanisms. Thus, this opens the door to go beyond the use of ad-hoc AI-augmentation techniques as is the case in today’s data-driven and information-driven networks. This, in turn, poses a fundamental question: “*How can we create reasoning-driven AI-native systems?*”. Our proposed reasoning-driven AI-native systems can thereby use less data, and more knowledge, in order to perform the various functions of the E2E wireless system.

One key challenge that must be addressed when answering the aforementioned question is the need to fundamentally transform the way in which data is viewed, processed, transmitted, recovered, and exploited at the level of the a network’s transmitter in receiver. In other words, creating reasoning-driven AI wireless systems must challenge the classical assumption that the wireless network’s transmitter and receiver (even those designed with AI) are simple “bit pipes” that act as a simple conduits of data bits, without exploiting a knowledge base that is built on the structure, linkages, and relationships between multiple low-level data points. Indeed, somewhat remarkably, despite all the effort directed towards the AI-nativeness of future wireless networks, at the level of transmitter and receiver, we still rely on very classical message construction and recovery mechanisms. One could argue that recent AI-based transmitter and receiver designs (e.g., using autoencoders [12]) are efficient in jointly learning transmitter and receiver implementations as well as signal encodings without any prior knowledge. However, although such mechanisms adopt deep learning to learn the transmitter, channel, and receiver, these learned building blocks fundamentally still perform the same classical information transmission tasks. In other words, the message construction and recovery mechanisms are “*learned*” to combat the channel’s uncertainty. Moreover, despite advances in AI, the operational functionality of communication systems

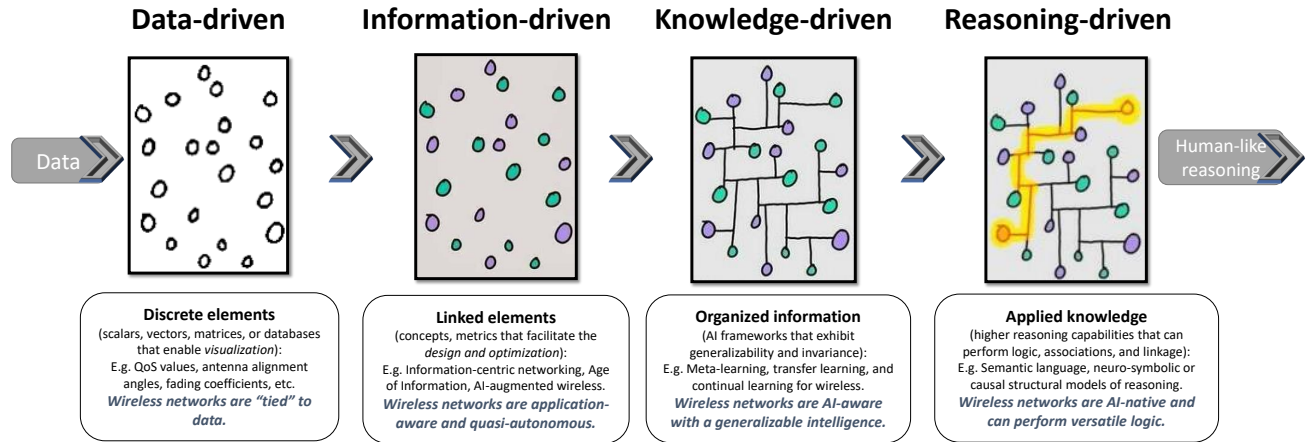


Fig. 1: Illustrative figure showcasing the evolution of wireless networks from data-driven ones towards reasoning-based ones.

remains tied to data, dependent on classical artificial neural networks (ANNs), with weak generalizability and reasoning abilities. Such frameworks fail to *organize information by characterizing the causal and associational characteristics of data* (See Fig. 1), thus, preventing the network from gaining knowledge that can be exploited in the future to operate with less input data. As a result, such *knowledge-agnostic* AI-native wireless network design approaches have a limited evolutionary potential and must re-engineered to mimic human reasoning, if we are to make a fundamental leap in future wireless technologies (6G and beyond).

Next, we discuss why the transmitter, receiver, and air-interface have retained their classical Shannon functionality to this date. Then, why the path towards AI-native, reasoning-driven networks requires a major leap from traditional communications to semantic communications, while leveraging the next wave of AI frameworks.

#### A. Why now? Why have we continued to rely on traditional communications so far?

Since its inception, the digital communication system problem, as posed by Shannon, has been a *reconstruction problem* because of the paucity in computing capabilities needed for more intelligent AI-guided tasks. That is, the fundamental goal of communications has hitherto been viewed as the capability of reproducing at one point either exactly or approximately a message transmitted from another point. Thus, the transmitter and the receiver have traditionally been designed in a fashion that relies solely on compression, transmission, and decompression. Moreover, the techniques used to encode the message only characterize the stochasticity stemming from the source, channel, and destination. For instance, the encoding performed at the physical layer characterizes the stochasticity at the transmitter, meanwhile,

such encoding does not represent the characteristics of the message conveyed, nor its context, i.e., the encoding does not contain any information related to the significance or meaning of the message. When augmented with AI, such approaches remain confined to data-driven designs.

In this conventional setting, one may not be able to efficiently convey the desired meaning of the transmitted messages (at the levels of the transmitter and receiver), thus leading to multiple repercussions on the overall end-to-end (E2E) communication system design. *First*, in many use-cases continual repetitive, back-and-forth transmissions are needed to transfer the necessary application information. That is, in the standard communication setting, the transmitter does not attempt to *automate the generation* of the message at the receiver. Also, the transmitter and receiver do not typically leverage their *memories*, i.e., the history of previous message transmissions/observations or observed past patterns in the data. This knowledge, if exploited properly, i.e., if the structure within historical message transmissions is *learned*, then the receiver can generate such a structure rather than *continuously reconstruct it* (see knowledge-driven case in Fig. 1). *Second*, in the current communication infrastructure, the receiver is, in general, *passive* and the communication mode is *asymmetrical*. In other words, the transmitter is typically in full control of the message generation and manipulation, and thus of its properties.

This passive behavior of the receiver and the asymmetry of communication links make the receiver susceptible to adverse channel conditions and any erroneous hardware or air-interface impediments. Augmenting this communication with a sense of symmetry, whereby the receiver can *learn* the structure of the messages and can leverage the history and context of the previously received messages, can potentially improve the robustness of communication with respect to channel and network irregularities. Broadly speaking, if the receiver is

endowed with the ability to generate its own messages and make its own conclusions, the E2E wireless network becomes less reliant on and susceptible to the wireless channel and its impediments. To exploit such history and learn context, the concept of *semantic communications* [13] can be leveraged. Semantic communication is a communication approach that promises to transform radio nodes into intelligent agents that can extract underlying semantics (meaning) in a datastream. That is, when communicating information, radio nodes leverage their reasoning faculties to identify the underlying structure of the message, and the role it plays in their knowledge base. Such a new form of communication can be very beneficial for scenarios in which the reliability of the link is intermittent. Examples of such scenarios include non-terrestrial networks (NTNs) (whose links are unreliable due to dynamic obstacles in space) and extremely high frequency (EHF) (millimeter wave (mmWave) and terahertz (THz)) whose links are highly susceptible to blockage and the radio environment in general. Here, semantic communications can overcome the intermittent behavior with the generative capabilities (via computing) of the receiver. For instance, under a semantic communication paradigm, radio nodes can take advantage of the concept of *semantic showers* to minimize back-and-forth communication with continual and reliable link (See Section II-B3 for more details). *Finally*, when messages are not perceived as a mere bit-pipeline, various types of context-related information can be discovered and exploited. Here, context is exploited at the level of the *data* itself (and its features) rather than at the level of the *application*. Indeed, even though the information at the application level can enable intelligent decision making at the radio node, such intelligence remains insufficient. That is, performing machine reasoning on the low-level bit-wise data opens the door for can enable the transmitter and receiver to discover the *causal roots of specific events in the messages*, regardless of the application-level information. Such causal information is a new input and knowledge that can be leveraged to steer the E2E communication system to attain particular goals or to simply automate the E2E wireless network operation. For example, if two robots are communicating with each other collaboratively, understanding the root cause of the messages sent from Robot 1 (one of the collaborative robots), can enable both robots to reach their ultimate goal more efficiently.

Clearly, it is desirable to transform today's communication systems to reasoning-driven semantic communication systems that intrinsically attribute meaning to the exchanged messages. In contrast to traditional communication systems that are driven by dynamic and uncertain communication resources, semantic systems leverage the computing resources and are founded, as will be evident from the rest of this paper, on the concepts of languages, reasoning, and causality. This transformation has the potential to substantially improve the efficiency and intelligence of future wireless networks, and ultimately achieve *more with less* via the convergence of the computing and communication resources. In the following

section, we will delve into the details of this transformation.

### B. From Transmitter/Receiver to Teacher/Apprentice

The essence of semantic communication is to humanize the communication between a transmitter and receiver so as to mimic knowledge-driven human conversations, interactions, and discussions. At a high level, semantic communication systems are ones that can perceive *the significance or the meaning* contained in a particular message. Semantic communication requires a rethinking of the communication problem with respect to the three levels introduced by Weaver [14]. In particular, right after Shannon proposed information theory, Weaver posited that communication involves problems at three levels [14]: i) *Level A*: The technical problem which measures the accuracy of the symbols of communication to be transmitted; ii) *Level B*: The semantic problem which concerns itself with the precision of the transmitted symbols with respect to the desired meaning; and iii) *Level C*: The effectiveness problem which measures the effectiveness of the received meaning on the conduct of the overall system. Traditional communication systems operated solely within the confines of Level A. However, if properly designed, semantic communication system can take advantage of advances in AI and computing power, and thus, potentially create communication systems that not only encompass all three levels proposed by Weaver, but go beyond them, via a *reasoning plane* (see Section VII), as more with less can be achieved. To do so, it is necessary to transform today's transmitter and receiver pair into what we propose to designate as *teacher and apprentice* nodes, whose capabilities are, the following:

- 1) *From a bit-driven transmitter to a knowledge-driven teacher*: The transmitter must be transformed from a bit pipe into a *teacher* capable of, first, disentangling multiple *semantic content elements* within to source data, i.e., separating different meaning, i.e., semantics, contained within a message. Then, for every semantic content element identified, the teacher must craft a *semantic representation* with desirable properties. Essentially, *the semantic content is the "meaningful" part of the data, and the semantic representation is the "minimal way to represent this meaning"*. This is similar to the way human beings try to find suitable words to describe their observations and ideas. Also, different semantic content elements could map to different modalities in the data. For instance, when hearing an audio recording of someone's voice, the tone of the voice can be one semantic content element, while the words pronounced are another semantic content element. A human being can easily disentangle and separate those two, and they can also understand the meaning of the words pronounced. Remarkably, today's communication system transmitter cannot identify or separate any underlying semantic structure or modality. It is thus desirable to re-engineer the transmitter to mimic a human being's reasoning capabilities (to the extent possi-

## Key Pillars of a Semantic Communication System

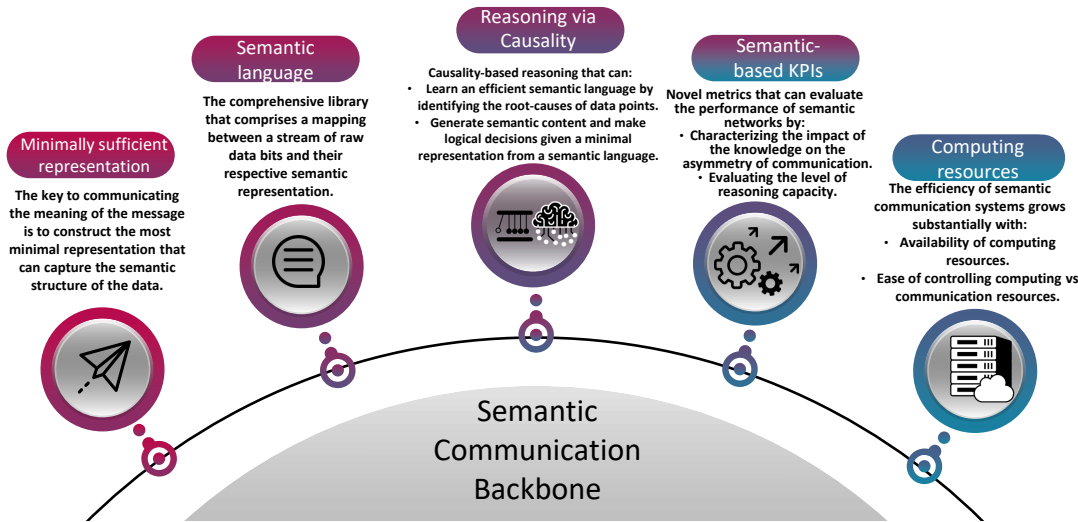


Fig. 2: Illustrative figure that showcases the key pillars of semantic communication systems.

ble). In this regard, on the transmission end there is a need for *reasoning*<sup>1</sup>: The facility that allows the transmitting agent to identify a semantic content element, distinguish it from others existing in the data, and devise an efficient representation of each of those identified contents. This is in stark contrast to the classical transmitter of today's networks that treats its input as a purely random and uncertain string of information and transmits it as a bit-pipeline that characterizes this uncertainty.

- 2) *From a bit-driven receiver to a knowledge-driven apprentice*: Similarly, at the receiver end, *reasoning* capabilities can transform the receiver into an apprentice capable of *understanding the minimal semantic representation* used by the teacher, i.e., mapping it to a semantic content element. Furthermore, the apprentice must be able to generate, via their computing resources, the semantic content element that results from the communicated semantic representation with the highest fidelity possible, e.g. if one of the semantic content elements were part of a hologram transmitted, the apprentice must be able to generate such a hologram with the same resolution that the teacher transmitted it. Moreover, as a result of the developed reasoning capabilities used to understand a semantic representation, the apprentice can use causal and associational (statistical) logic to perform various projections and decisions across the networking stack. Such causal and associational logic is inferred from the progressively built knowledge base and exerted on the received semantic representation.

<sup>1</sup>It is important to note that this is a broad definition for reasoning. In Section IV, we more concretely elucidate the definition of reasoning from a causal perspective.

- 3) *From a bit-pipeline to a semantic language*: In semantic communications, the smallest distinct meaningful element is a semantic representation. Moreover, a series of representations constitutes a *semantic language*. Semantic languages will mimic natural languages but they should be less focused to syntax and pragmatics in order to automate processes better (see Section IV-C). Moreover, the semantic representations of a semantic communication language must satisfy three key properties:

- a) **Minimalism**: The capability of characterizing the structure found in the information with the least number of language elements possible (and their equivalent bits). This characterization must be performed in a way to reduce the number of exchanged messages in the long run as well.
- b) **Generalizability**: The capability of representing a particular underlying structure (or understanding one at the receiving end) while being invariant to changes in: a) distribution, b) domain, and c) context. Notably, here, context can be viewed as the theme encapsulating various semantics that share a common denominator (for example, the context of messages received by a robot can be a set of steering actions on a tennis court). Hence, generalizability means that, when a radio node has learned and established a particular representation  $Z_i$  for a semantic content element  $Y_i$ , and it can then use such a mature and consistent representation to describe this semantic content element irrespective of the distribution, domain, or context it is extracted from. Hence, the node is now able to generalize its knowledge across multiple, previously unseen and unknown domains, distributions, and contexts. This

## Semantic Communication Systems: Effectiveness/Efficacy vs. Efficiency

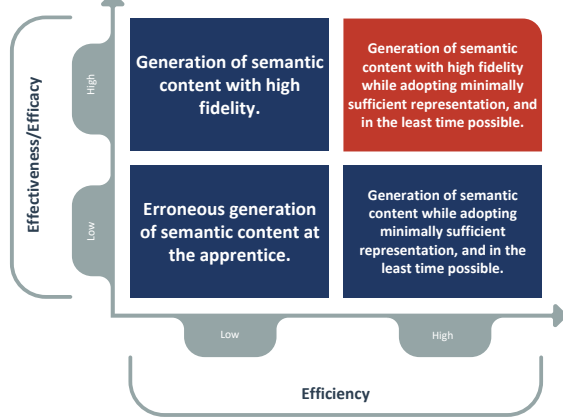


Fig. 3: Illustrative figure that showcases the differences between efficiency and efficacy/effectiveness.

mimics the behavior of words in a natural language that can *generalize* and apply their knowledge base to describe any occurring event, even if previously unobserved. For instance, a human being can identify and describe a “loud voice”, regardless of what the voice is saying or the environment it is observed in. This generalizability capability represents the epitome of reasoning in regards to making logical conclusions.

- c) **Efficiency:** The ability of the apprentice to re-generate the information with high fidelity, in the least time possible. In other words, the resolution of the data generated at the apprentice must be equal (or better) to that which could be recovered by a classical receiver. For instance, if the apprentice is trying to re-generate an audio clip, the resolution of the audio clip must be at the same resolution intended by the teacher. In other words, given that a *reasoning* process has been added, there should not be a tradeoff between the quality of the content recovered and the minimalism achieved via semantic communications. It is important to note, that while the effects of efficiency can be measured via concrete metrics (which are defined via the semantic impact in Section VI), it is more difficult to evaluate the efficacy of the system. That said, for the sake of brevity, in this tutorial when discussing “efficiency”, the term is considered to mean *efficacy and efficiency* simultaneously, i.e., the upper right block of Fig. 3.

As a result of the aforementioned desirable properties, a well-designed semantic language can achieve *more with less*.

A summary of the key pillars of semantic communication systems are shown in Fig. 2. Essentially, a classical communication system revolves around the channel and the communication resources, which are dynamic and governed by uncertainty. In contrast, the major thrust of a semantic communication network is its reasoning (and knowledge) capabilities that can be achieved via causality. In essence,

as shown in Fig. 2, on the one hand, a reasoning radio node will communicate via a semantic language. A semantic language is a comprehensive library that maps every semantic content element in a raw datastream to a *minimally sufficient representation*. A minimal sufficient representation is needed to achieve minimalism both on the short and long term. Ideally, such a representation would: a) minimize the communication resources on the short term, and b) enable the apprentice’s reasoning capabilities to scale up, to ultimately generalize and perform versatile logic operations. On the other hand, properly deploying reasoning capabilities is at the helm of evaluating the semantic communication network via a suite of novel semantic-based key performance indicators (KPIs) which can properly capture the new reasoning dynamics of the performance. Furthermore, causal and associational logic via reasoning is not *feasible* without abundant computing resources that exhibit a flexibility in control, in contrast to communication resources.

In a nutshell, we have thus far elucidated the roles of the teacher and apprentice in a semantic communication system. We have also highlighted the need for a semantic language and overviewed its unique characteristics. Next, we highlight the main contributions of this article.

### C. Contributions

The main contribution of this article is a novel and holistic vision that articulates fundamental principles necessary to build next-generation reasoning-driven, AI-native semantic communication networks. In particular, we *first*, investigate the key tenets necessary to extend today’s classical information theory towards a semantic information theory. This extension is performed via a migration from today’s *bit-pipeline* to a *semantic language*. Second, we scrutinize the reasoning foundations that are imperative for the communication of a semantic language. These foundations are centered around the migration from *data-driven* networks towards *knowledge-driven and reasoning-driven* ones. In essence, organizing *information* is majorly influenced by the capability of a radio node to unravel causal and associational relationships and logic. In this regard, we propose rigorous reasoning techniques that must be adopted in gradually building a language, to ultimately reach a *generalizable, minimal, and efficient* semantic language. Here, we particularly shed light on the significance of causality to rise up in the *reasoning ladder* (see Fig. 11). Third, we propose a suite of novel *semantic KPIs* for evaluating the performance of AI-native, reasoning-driven systems, and optimizing future semantic networks. Finally, we discuss how one can build scalable semantic communication networks while bringing forth novel approaches and concepts to address several computing, control, and networking challenges. Furthermore, through our contributions, we answer the following fundamental questions:

- *How do we extend classical information theory to capture semantic information?*

The performance of today’s communication systems is

evaluated based on principles derived from Shannon’s information theory. That said, information theory is built on the premise of defining “information” as a mere “uncertainty” that does not perceive meaning or structure [14]. In consequence, we investigate and characterize the equivalents of today’s “information” and “entropy” in a semantic communication system. Then, we discuss how these novel concepts modify the way communication is intrinsically viewed and evaluated.

- *Why do we need a semantic language? How is it different from a natural language?*

For a radio node to become capable of communicating using a semantic language, it must be able to: a) extract semantic content elements from the data, b) map into a minimal semantic representation, and c) understand a semantic representation occurring in various domains, contexts, and stemming from different distributions. We show that a semantic language is fundamentally different from a natural language. The atomic unit of a semantic language is a *representation* that captures *the structure and variability of the represented semantic content element*. Meanwhile, the atomic unit of natural language is a word. Limiting a semantic language to a natural one, would constrain it in syntax and wording which (unlike causal and associational logic) are governed by deterministic rules.

- *How do we semantically process data, and how can we build a semantic language?*

Many tools can enable extracting a semantic representation from data. However, the right approaches must be able to create a semantic representation that is minimal, efficient, and generalizable. To obtain such a representation, one must be able to characterize the causal and statistical properties of the data. Thus, after surveying the set of existing tools for representing semantic information, we expose the fundamentals of causal representation learning and its accompanying benefits, challenges, and future directions for building next-generation semantic communication networks.

- *How do we move from data-driven intelligence towards knowledge-driven reasoning?*

Moving from data-driven intelligence towards knowledge-driven reasoning requires engineering the semantic language based on a model that can characterize causal and associational logic. As a result, we demonstrate that mapping a language to a *structural causal model (SCM)*, enables exploiting the concepts of *interventions and counterfactuals* from causal logic. Such concepts allow building a semantic control plane, whereby instead of classical acknowledgements and non-acknowledgements, the apprentice can gather information about the structure of the previously conveyed representation. This process enables a gradual acquisition of a language at the apprentice, which ultimately leads to elevating the radio node in the causal reasoning ladder.

- *How do we evaluate the performance of semantic communication systems?*

The evaluation metrics of classical communication systems have heavily relied on Shannon’s information theory, however, future evaluation schemes for semantic-based systems must capture the structure of the semantic representations and the reasoning capability of the teacher and apprentice. Thus, we propose three novel semantic-based metrics that enable characterizing the semantic impact a particular representation can generate (which characterizes the gain in time and resources when relying on a semantic representation versus classical data), the communication symmetry index, and the reasoning capacity of a semantic communication link.

- *How do we scale semantic communication systems to current and future large-scale cellular communication networks?*

Today’s 5G cellular networks are characterized with a separated control and user plane. With the introduction of semantic communications, many fundamental changes are necessary ranging from the need to integrate causality and reasoning to designing an expressive yet minimal semantic language. One key change is the need to introduce a novel reasoning plane, that would be sandwiched between the control and the user plane. Based on the real-time inference performed in the reasoning plane, the control plane is fed with information that enable radio nodes exchanging interventions and counterfactuals. These are queries that replace acknowledgements and non-acknowledgements and enable the teacher and apprentice to build and learn a language. Also, with the introduction of semantic communication, today’s open-radio access network (O-RAN) concept will evolve further to account for real-time, near-real-time, and non-real time intelligence.

#### D. Prior Works

Recently, a number of surveys and tutorials related to the concept of semantic communications have appeared in [15]–[22]. The authors in [15], presented a view on semantic communications within three communication modalities: human-to-human, machine-to-human, and machine-to-machine. The authors in [16] present the developments of deep learning (DL)-enabled semantic communications for multi-modal data transmission, including text, image, and audio. In [17], the authors overview a semantic signal processing framework that can be tailored for specific applications and goals. In [18] and [19], semantic and goal-oriented communications were overviewed while highlighting the network benefits in terms of reliability and effectiveness. The work in [20] presents key methods for performing feature extraction based on semantic communications. In [21], the authors analyze semantic communications from an information-theoretic perspective. The authors in [22] discuss how 6G technologies can drive the development of semantic communications.

While the works in [15]–[22] are interesting they have

## Six Essential Layers for Building Future Semantic Communication Systems: Foundations, Approaches, And Scalability

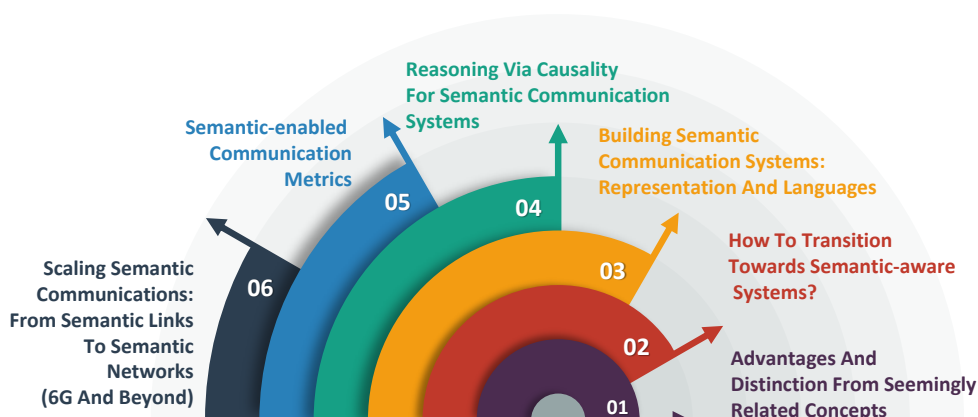


Fig. 4: Main sections of this tutorial via the 6 Essential Layers for Building Future Semantic Communication Systems.

not considered various concepts and fundamentals that are necessary to define, design, and ultimately deploy semantic communications systems:

- **Representation:** While the works in [15]–[22] acknowledge the need to depart from a latent bit-pipeline, such works do not lay the language pre-processing techniques that are needed to move from entangled raw datastreams to *learnable* datastreams that can be used to learn a semantic language. Moreover, the previous works in [15]–[22] fail to articulate the necessary measures that enable the apprentice to *understand* a representation, i.e., leverage it to efficiently generate semantic content elements via their computing resources.
- **Semantic Language:** While some of the works in [15]–[22] acknowledge the need for a language, these works fail to explicitly define a semantic language, its corresponding characteristics, and how it can improve the efficiency of communications. Moreover, to mimic human conversations such works often confuse natural languages with semantic languages. In contrast, this tutorial will be the first to extend Shannon’s information theory to the semantic communication domain by investigating the necessary measures to gradually build a semantic language. In fact, the main goal of a semantic language is to express a minimal semantic representation vis-à-vis the raw datastream and its contained semantic content element. Such a language is characterized with *minimalism*, *generalizability*, and *efficiency*. Also, our work will be the first to further elucidate the fundamental tenets of this language by distinguishing from natural languages with respect to syntax, pragmatics, and semantics.
- **Reasoning and Causality:** The prior art [15]–[22] does not provide any concrete technical approaches to perform

reasoning. In particular, they do not put into perspective the importance of *causality* in the data. In many ways, existing works limit themselves to statistical and associational relationship in the data that fail to unravel the underlying structure of the data. Also, relying on such statistical relationships may lead to spurious representations. Instead, in this tutorial, we are the first leverage the concept of causality. In essence, ultimately enables communication nodes to recognize the root causes of specific datastreams via the notions of interventions and counterfactuals.

- **Semantic KPIs:** All the prior works in this area in [15]–[22] still rely on classical KPIs such as rate, reliability, and latency. These KPIs cannot capture the reasoning capability of radio nodes, nor can characterize the level of communication symmetry. In contrast, we propose a suite of semantic-based evaluation metrics that enable characterizing the performance bounds of any semantic communication system. For example, we derive a metric called “reasoning capacity” that can ultimately have an impact that is higher than Shannon’s capacity because its reliance on computing resources (in contrast to communication resources).
- **Scalable semantic communications:** The works in [15]–[22] are limited to one source and destination, meanwhile the potential of semantic communications cannot be fully unleashed unless it is considered over large-scale cellular networks. As such, in this work we elucidate the challenges and the opportunities regarding scaling semantic communications over a wireless network.

With the currently established literature, designing and building future semantic communication systems from the ground up is an extremely strenuous and difficult task. Essentially, there is



a major lack in existing works that investigate novel knowledge and reasoning frameworks, which constitute a central pillar to propose a comprehensive framework for semantic communications. It is thus necessary to investigate the overarching measures needed to successfully usher the birth of semantic communication systems in beyond 6G systems. Next, we examine the preliminaries of semantic communications by highlighting its underlying benefits and distinguishing the concept of semantic communications from alternative concepts that have recently emerged.

### E. Organization

The rest of this paper is organized as shown in Fig. 4 and as follows. In Section II we discuss some of the advantages of semantic communications and its relationship to existing techniques. Then, in Section III we discuss novel concepts and definitions that enable a smooth transition from classical communication systems towards semantic communications system. Subsequently, in Section IV we thoroughly investigate novel views and techniques that enable establishing an expressive semantic representation and language via scrutinizing the structure in the data. Then, in Section V, we demonstrate how to leverage, for the first time, the concept of causality to ultimately equip radio nodes with a reasoning capability. Furthermore, in Section VI, we propose a suite of novel semantic metrics that enable evaluating emerging semantic communication system. Then, in Section VII we develop some of the key techniques that enable the design and deployment of semantic communication networks at scale. Finally, conclusions and recommendations are drawn in Section VIII.

## II. SEMANTIC COMMUNICATIONS: ADVANTAGES AND DISTINCTION FROM SEEMINGLY RELATED CONCEPTS

Before delving into the main technical components of semantic communications, we first distinguish the concept of semantic communications from alternative frameworks and concepts that were recently proposed for next-generation wireless systems. Then, we overview the benefits of semantic communication networks.

### A. What is NOT Semantic Communications?

At first glance, semantic communications can seem like an incremental variant of known approaches and techniques. In this subsection, we attempt to demystify this confusion by highlighting the fundamental differences between such techniques and semantic communications. In Fig. 5, we summarize our answer to the questions of what is and what is not semantic communications.

1) *Semantic communications is not data compression:* In a classical communication setting, according to information theory [23], the process of data compression (also known as source coding or bit reduction) is the process of encoding information<sup>2</sup> using fewer bits than the original datastream

<sup>2</sup>Here information denotes the Shannon definition of information, and thus is designating the concept of uncertainty which will be elaborated in Section III-B1.

representation. This process is performed by exploiting the statistical redundancy in the data bits to ultimately represent data without any loss of information. As such, the process becomes reversible at the receiver. While data compression shares some common ground with semantic communications with respect to minimalism, i.e., minimizing the size of data transmitted, both concepts are fundamentally different:

- Data compression achieves minimalism, i.e., shrinking the size of a particular datastream, by identifying and eliminating statistical redundancy. For instance, the most prominent lossless compressors employ probabilistic models such as prediction by partial matching [24]. Notably, there is a close connection between data compression and ML, in that they both specifically attempt to predict the posterior probabilities of a sequence given its history. Therefore, a duality arises between data compression and ML, to the extent that some works in [25], consider data compression as a key method that can be used in ML for tasks like clustering and classification. That said, as a concept, data compression does not bear any learning ability that contributes to a particular training memory or a trained model. In other words, from an ML perspective, data compression techniques often intend to *overfit* since their only goal is to shrink the current datastream and not the future ones. Thus, while data compression can be a component within an AI technique, it does not exhibit *learning* characteristics, let alone reasoning. From a minimalist perspective, data compression could even be more beneficial than semantic communications on the short term. However, it is unable to instill contextual information and knowledge-driven memory on the receiver. In many ways, on its own, data compression is restricted to the realms of data-driven and information-driven networks from Fig. 1.
- In contrast to data compression, in semantic communications, instead of identifying data redundancies and compressing them, patterns that map to structure and semantic content are identified, learned, and then represented with a semantic representation. In essence, in source coding and data compression, the goal is to *overfit* to the statistical characteristics of the datastreams. Meanwhile, semantic communication's ultimate goal is to characterize structure, and, thus, the focus is not on the pure randomness exhibited in the data. In fact, these random data points are better transmitted classically as we explain in Section IV. Furthermore, semantic communication achieves "minimalism" as a byproduct of the semantic representations transmitted which: a) Comprise a fewer number of bits in total, b) Serve to teach the apprentice to learn, generate, and ultimately automate the task or message at the receiver. Consequently, the characteristics of the adopted representation and the acquired reasoning capabilities enable minimalism via: a) Minimizing the number of bits per transmission, and b) Minimizing the total number of transmissions necessary to convey a message or achieve

## What is (is not) Semantic Communication (SC) Systems?

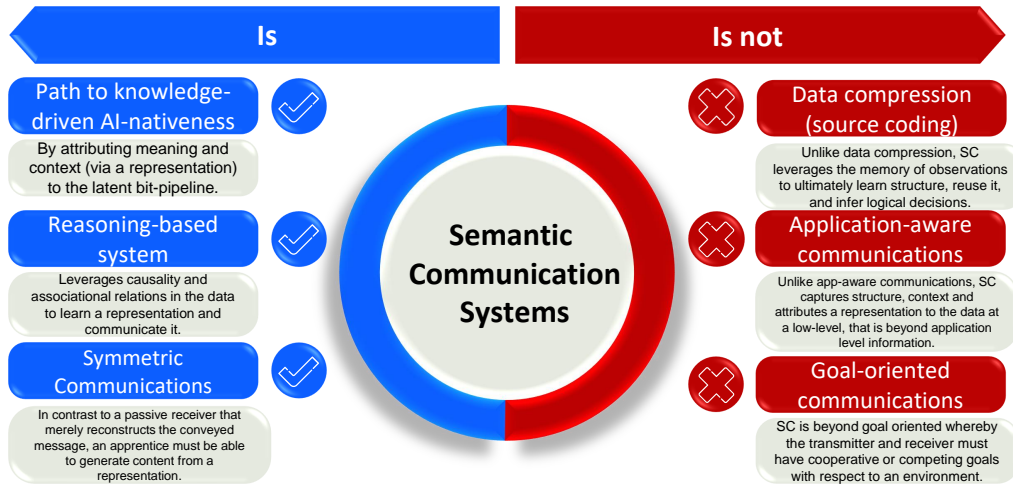


Fig. 5: Illustrative figure showcasing what is (is not) semantic communications.

a task. Hence, semantic communication networks achieve minimalism via different mechanisms that go beyond a compression of the number of bits in a packet. Finally, when radio nodes operate based on organized knowledge, such nodes can make more informed logical conclusions across the networking stack, a feature only possible with reasoning-driven semantic communication networks.

2) *Semantic communications is not only an “AI for wireless” concept*: AI has been used for many wireless-related problems in the past few years. This includes AI and ML for channel estimation, beamforming, network management, receiver design, etc. In essence, in all of these tasks, the design, performance, or optimization of the task was procured via an AI or an ML tool instead of traditional numerical methods. AI-enabling such tasks was shown to improve the accuracy, precision, or relevant KPI, however, the fundamental functionality and dynamics of the corresponding wireless task was kept the same. For example, AI-enabled channel estimation is fundamentally an improved approach to perform classical channel estimation solutions (e.g. as performed in [2]), yet the core task to be performed remains the same. In contrast, semantic communication systems will not be an additional AI-enabled layer on top of an originally existent task. In other words, semantic communications is not an AI-improved air interface.

With semantic communications, user equipments (UEs) and base stations (BSs) do not need to rely continuously on the channel. For instance, for scenarios in which the radio nodes have established a solid knowledge base, the need for channel estimation becomes minimal. Moreover, instead of relying on bit-driven signaling messages to sense the channel, a continual understanding of the context of previous messages

can potentially establish an awareness of the physical environment. Such awareness can be leveraged to *learn* the channel characteristics while relying on computing resources. In essence, the *mechanism of communication tasks fundamentally changes with the introduction of semantic communications* (as previously explained for channel estimation as an example). This can be observed via the following insights:

- In contrast to a classical AI-based transmitter that still relays bits “as is”, the teacher in a semantic communication network can now attribute a semantic language to the raw bit-pipeline previously observed. Meanwhile, an AI-layer can only add a prediction margin to the classical tasks of compression, transmission, and reconstruction.
- When communication relies on a semantic language, context becomes of relevance. As such, the more consistent the theme of the conversation is, and hence the context, the larger the improvement in the reasoning capability at the source and destination. In contrast, classical AI-augmentations alone are not aware of the concept of context, i.e., the logical theme surrounding the data structures learned by exploiting knowledge base. Instead, data/information-driven AI algorithms are ultimately reliant on their data input and their corresponding statistical properties.
- Classical AI for wireless attempts to improve the performance after being trained a priori via large datasets, or via a trial-and-error phase (e.g. RL) over the inputs from a specific environment (e.g. channel, spectrum, QoS values). Instead, semantic communications *gradually* builds a language between the teacher and the apprentice. This gradual language construction enables the teacher and the apprentice to organize and build their knowledge base.

Thus, the radio nodes now acquire *human-like* reasoning faculties. That is, a radio node can now: a) make conclusions according to its knowledge base (not data), and b) communicate its needs based on such conclusions.

Clearly, based on the above key observations one can conclude that semantic communications is beyond a simple use of an AI algorithm for a wireless task. Given that semantic communications enables radio nodes to build a knowledge base and communicate a language, the mechanism of communication fundamentally improves.

3) *Semantic communications is not only goal-oriented communications*: A goal oriented communication system involves a number of agents that interact and exchange messages to achieve a joint goal or separate goals that include the same environment. For example, two robots can interact with each other to execute a common mission. Here, in contrast to sending the information gathered by sensors bit-by-bit, the robots can exchange multiple feedback messages of their current semantic action, their next expected outcome, all while achieving a unique joint goal. In a goal-oriented framework, the nodes, e.g., the teacher and apprentice can also be achieving two separate goals. Much of the early-on work on semantic communication has equated it with such goal-oriented communication systems [26]–[29]. However, there are fundamental differences between the two concepts. In some sense, goal-oriented communications falls under the umbrella of semantic communications. For instance, in every goal-oriented communication system, the nodes will have to embed semantic representations to ultimately achieve a particular goal. In contrast, under the broader auspices of a semantic communication system, the generation and communication of semantic representations is not necessarily done for the purpose of serving a system-wide goal

In this regard, limiting the concept of semantic communications to the confines of goal-oriented systems will therefore unnecessarily limit its use to a subset of use-cases that have a competitive or cooperative nature. Meanwhile, there are many instances in which the teacher and the apprentice do not necessarily share any joint goals nor interact with a common environment. For instance, the teacher can be a server that is transmitting highly-data intensive content (e.g. extended reality (XR) content) to a particular user. Here, every standalone content transmitted can have an entirely different goal, and there are no cooperative or competing goals between the teacher and the apprentice. Yet, in this case, semantic communications can still be used to: a) rely less on the channel to transfer massive information content, b) empower radio nodes with reasoning to make versatile decisions, which can enhance network's capability in meeting the stringent requirements of future applications.

4) *Semantic communications is not only application-aware communications*: Implementing the context of information within the transmission of messages may seem at first glance similar to the traditional concept of application-aware communication. In fact, there are many prior works (e.g. [30]–

[35]) that have fine-tuned the network optimization process to address application-level requirements. For example, in XR applications, the XR content transmitted by users may exhibit a particular correlation. Here, some works such as [31] exploit this correlation to ultimately improve the management of uplink and downlink wireless transmissions. Notably, it is important to distinguish between the “context-awareness” concept defined by such frameworks and the one granted with semantic communications systems, as the former is a mere application and use-case specific awareness. In contrast, in semantic communication systems, “context” is a concept defined with respect to the low-level structure of exchanged datastreams between the transmitter and the receiver. Such low-level intelligence opens the door for an inter-application, intra-application, and out-of-domain generalizability. In other words, a radio node can leverage the meaning attributed to low-level data corresponding to service *A* by using it to improve the E2E performance for service *B*.

On top of gaining generalizability, the “awareness” gained from application requirements, as done in classical application-aware works [30]–[35], majorly relies on statistical learning and cannot be easily extended to solve multiple challenges across the open systems interconnection (OSI) model. Meanwhile, semantic communications intrinsically relies on causal and statistical relationships in the context of the data. For example, after acquiring the structure of the data, the radio node can track the root cause of a mismanaged resource orchestration to ultimately predict and proactively prevent a beamforming error from happening. The radio node can also conclude information about the type of terrain in which communication is taking place (e.g. rural, urban).

Furthermore, some works [36]–[38] view semantic communications through a *significance perspective*. Such works [36]–[38], consider metrics like AoI, value-of-information, and other time-oriented metrics to be indicative of *significance* of the use-case. While such metrics are useful in enhancing the performance of time-critical communications, they are still limited to certain use cases, and they do not possess the generalizability and reasoning capabilities needed for semantic communications. That is, AoI is a networking metric that does not unravel low-layer information (which contains the structure of the data). Such networking metric alone, do not allow a radio node to build a knowledge base, thus such radio nodes fail to perform any reasoning-driven tasks. Additionally, such metrics are still highly dependent on communication resources and are mainly constrained to the framework of particular use-cases and time-critical communications. Thus, for instance, such significance measures cannot universally enhance the performance irrespective of the use-case or its time-criticality. In many ways, AoI and its variants still lie in the scope of information-driven networks from Fig. 1. Clearly, semantic communications equips nodes with an intelligence that has breadth and depth that is beyond the one gained with application-aware communications.

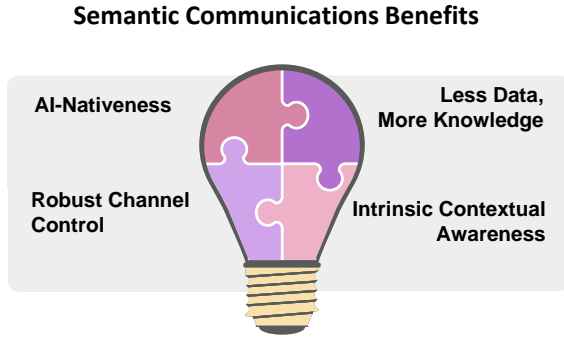


Fig. 6: Illustrative figure showcasing the benefits of semantic communications.

### B. Benefits of Semantic Communications

Semantic communications will bring forth many key benefits for future communication systems, in what follows we detail those benefits.

1) *Achieving Reasoning-Driven AI-Nativeness*: 6G and beyond systems must be AI-native across their protocol stack, in the sense that every single component, layer, and structure in the network must be designed, deployed, and optimized via AI. Nonetheless, today’s industry inaccurately describes future AI-native systems, as ones reliant on data or information, as shown in Fig. 1. Such classical AI approaches are rigid, big-data dependent, and are knowledge-agnostic. Meanwhile, as we discussed in Section I we define AI-native networks, as ones that *rely on less data but more knowledge*. That is, the overarching goal is to create *reasoning-driven* systems as shown in Fig. 1. Effectively, semantic communication can provide a path to a next-generation of knowledge-driven and reasoning-driven radio nodes. This can be done by: a) Transforming today’s bit-pipeline communication to one that relies on a semantic language. In essence, to communicate a semantic language, radio nodes must utilize the input semantic representation to *computationally generate* semantic content elements, and b) Performing reasoning which is attained by extracting/understanding semantic representations. This enables radio nodes to capitalize the context of the semantic language, as well as their knowledge base to make versatile decisions across the networking stack. Clearly, radio nodes that communicate a semantic language and are reasoning-driven, can potentially reach the real-time prediction, automation, and agility needed for 6G and beyond systems.

2) *Intrinsic Contextual Awareness*: One key benefit of semantic communications is that it can grant radios with an intrinsic awareness of the contextual information of their data transmission. This means that radio nodes (e.g. BS or a UE) at any point in time, become aware of the context, i.e., the revolving communication theme of recurrent messages. For

instance, if a person is sending pictures of themselves and their pet in a park, the spatio-temporal settings, as well as the revolving picture colors have an underlying *family of structure*. This consistent family of data structures is the *context* of the messages transmitted in this scenario. An awareness of context is particularly important for Internet of Everything (IoE) services as they require the co-design and control of different functional modalities of communication, computing control, sensing, and tracking. Herein one can for example improve the optimization of control, localization, and sensing messages. The converse can also be true: Properly extracting the semantic information gained from sensing or tracking messages, can further improve the optimization of communication resources. Moreover, this contextual awareness enables the radio nodes to be self-governed and self-optimized while being offline, i.e., such radio nodes can rely on their knowledge base to make future logical decisions without a continuous connectivity. For instance, in many settings, the radio node can extract the spatio-temporal changes from the semantic content of a communication message. Hence, owing to the spatio-temporal changes causally learned, the network can minimize the number of control acknowledgements needed as well as frequent tracking and signals. This can ultimately improve the spectrum and energy efficiency of the sensors, radar BSs, and control BSs.

3) *Robust Channel Control*: To date, wireless communication systems have always been governed by the channel and its uncertainty. In essence, most of today’s wireless communications research efforts aim to analyze and optimize the network performance to ultimately render the information transfer task more efficient over the *uncontrollable channel*. Nonetheless, in semantic communications, the communication link is no longer an asymmetrical link that mainly relies on decoding every single bit to recover the identical message transmitted. In fact, as a byproduct of the convergence of computing and communication resources, the teacher and apprentice’s center of gravity will move more towards *controllable computing resources*. This shift towards computing enables communication systems to rely less on the classical 3GPP concepts of reliability and continuity. Thus, a higher level of independence and robustness to various channel conditions can be gained. Such independence and robustness is exhibited by two mechanisms: First, given that communication now depends on a semantic language, the teacher can offload *semantic showers* to the apprentice. Such semantic showers would contain the language that *explains* to the apprentice the meaning of the message. Then, by utilizing their computing resources, the apprentice would generate the elements of the service content. This minimizes the back-and-forth exchange of user and signaling messages needed to continuously convey information via a communication channel. In fact, such semantic showers can be particularly exploited in future networks that will rely on intermittent THz or mmWave. Also, such showers enable improving the intermittent service of NTN that are needed to *bridge the digital divide* and provide global coverage. Second,

Table I: Common lexicon used in semantic communication systems.

Vocabulary	Definition	Mathematical Expression
<b>Teacher</b>	A transmitter with reasoning capabilities. A teacher is capable of first disentangling multiple semantic content elements to be transmitted, i.e., separating different meaning contained within the message. Then, for every semantic content element identified, they will craft a semantic representation with desirable properties.	$b \in \mathcal{B}$
<b>Apprentice</b>	A receiver with reasoning capabilities. The apprentice can map the conveyed semantic representation to a semantic content, i.e., mapping the minimal representation to its corresponding meaning. Then, the apprentice can generate the content at the destination with the same fidelity initially produced at the source.	$d \in \mathcal{D}$ .
<b>Semantic Content Element</b>	The meaning of a specific datastream (or a label denoting this meaning).	$Y_i$
<b>Fidelity of Information</b>	A high fidelity corresponds to recovered information with an equivocal resolution of the data type and content transmitted, e.g. for image data, a high fidelity is an image recovered in its original resolution.	N/A
<b>Semantic Representation</b>	A representation that has desirable properties, and that is capable of “describing” the meaning of a datastream. This representation must be sufficient so that the apprentice can generate the semantic content without sacrificing the fidelity of information.	$Z_i$
<b>Semantic Language</b>	A dictionary (in terms of data structures) that maps every raw datastream to its corresponding semantic representation.	$\mathcal{L}$
<b>Semantic Didactics</b>	A combination of a stream of a semantic representations, complemented with a raw datastream sent by the teacher to gradually teach the apprentice the semantic language.	N/A
<b>Disentanglement</b>	The process of separating multiple semantic content elements in a single datastream.	Eq. 10
<b>Reasoning</b>	<ul style="list-style-type: none"> <li>• On the transmitter’s end: A reasoning-driven teacher capable of disentangling multiple semantic content elements within a datastream, and attributing each one a semantic representation</li> <li>• On the receiver’s end: A reasoning-driven apprentice, mapping the received semantic representation to a semantic content elements, and generating such content with high fidelity.</li> <li>• Reasoning at both teacher and apprentice must allow them to use their built knowledge base to perform logic operations and draw logical conclusions across the networking stack.</li> </ul>	Proposition 3
<b>Context</b>	A single context corresponds to a family of structures that is shared upon multiple recurrent datastreams. It can also be viewed as the theme encapsulating various semantics that share a common denominator.	N/A
<b>Dynamic Reasoning</b>	The capability to perform reasoning over varying context.	N/A
<b>Minimalism</b>	The capability of characterizing the structure found in the information with the least number of bits possible. This characterization must be performed in a way to reduce the number of exchanged messages on the long run	N/A
<b>Efficiency</b>	The ability of the apprentice to re-generate the information with high fidelity, in the least time possible. That is, the resolution of the data generated at the apprentice must be equal (or better) to the one that could be recovered by a classical receiver.	Large semantic impact, i.e., $\iota_\tau > 1$ (See Definition 12 and Proposition 2).
<b>Generalizability</b>	The capability of representing an underlying structure when dealing with datastreams of varying: a) distribution, b) domain, and c) context. This embodies the ability of a radio node to generalize and use its knowledge base to draw conclusions via a mature semantic language that can identify semantic content elements irrespective where they are drawn from.	Definition 10.

on top of minimizing the back-and-forth messaging, radio nodes in a semantic setting can leverage their knowledge base to *correct erroneous semantic representations*. Instead of solely relying on error correcting codes, when a semantic

representation results in a *spurious semantic content element*, i.e., inconsistent with the current context of communication, the radio node could deploy its logic to *predict* what representation the teacher was intending to send.

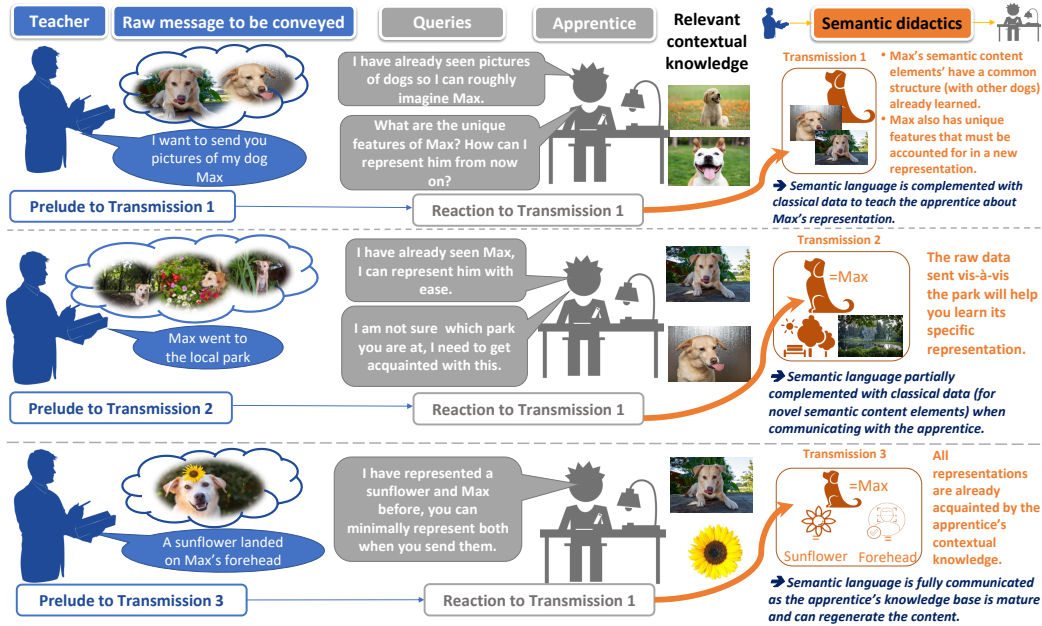


Fig. 7: Illustrative example of a reasoning-based communication framework that is gradually building a semantic language whereby they migrate from relying on discrete data elements to organized, linked, and logic inducing knowledge bases.

4) *Less Data, More Knowledge*: One key benefit and feature of the semantic language is *minimalism*. This minimalism is exhibited via two mechanisms: First, as will be demonstrated in the sequel of this work in Section IV, communicating via a language with tolerable complexity, yet significant structure imposes the *minimal sufficiency condition on representations*. In essence, a semantic representation is *the teacher's minimal description with regards to the meaning held by particular content element in the datastream*. That is, such a representation must be characterized via the *smallest number of bits*, while also being able *accurately describe the semantic content*. Second, given that radio nodes can leverage their knowledge base to perform *reasoning*, the apprentice can draw conclusions based on the context communicated. These conclusions can minimize the number of redundant back and forth messaging. For instance, in an autonomous vehicle setting, instead of tuning the car's direction in real-time and exhaust communications resources, the controller could send a representation that depicts the mode of driving in the future (e.g. drive straight for the next hour in high cautiousness). Here, the apprentice would leverage their computing resources to *generate* how driving straight and in high cautiousness is exhibited (these exhibits are the semantic content elements). Thus, a radio node is operating in a "less data, more knowledge", whereby the data stemming from representations and back-and-forth messaging is minimal, yet the knowledge base of a radio node is comprehensive and elevates the intellect of this node. Remarkably, this "less data" realm *minimizes the reliance on the spectrum*, and the need to open new spectrum bands in order to respond to the exponential data rate increases from one cellular generation to another. Thus, one long-term key benefit of semantic communications

is a minimization in the need for the more spectrum as well as complex dynamic spectrum sharing schemes (technical and regulatory) whenever new wireless technologies or use cases appear.

Next, we overview the fundamental measures that must be revisited in classical information theory to transition towards semantic communication systems.

### III. HOW TO TRANSITION TOWARDS SEMANTIC-AWARE SYSTEMS?

#### A. From bit transmissions to knowledge-driven human conversation

One caveat of a classical communication scenario is constrained to a bit agnostic representation of the message. For instance, if the transmitter were to send a photo of a dog to the receiver, then the photo must undergo signal processing steps to finally be represented on the bit and consequently packet level. Then, an erroneous reconstruction of the image at the receiver may result from any singular bit error. These bit errors can stem from errors at the transmitter, channel, or receiver. Also, such errors can have a hardware, software, or networking nature. In other words, a "bit-error" can result from a hardware/software defect or a networking bottleneck. This phenomenon can be called *data-blindness* at the transmitter, receiver, and air-interface levels. Moreover, given that the reconstruction process is oblivious to the context of the bits, e.g., whether a bit represents a dimension in the background or foreground; the reconstruction process is highly susceptible to these aforementioned bit errors. That is, current error correction schemes that attempt to minimize errors, lack an AI-foundation that unravels the root-cause of the errors/bottlenecks in the net-

work and attributes context to it. Hence, such error correction mechanisms are *memory-less* and can only minimize errors based on “current datastreams”. Thus, these schemes cannot leverage recurrent and semantic errors to ultimately improve the long-term system performance. In contrast, if bits were to become aware of semantics and context, the robustness to errors and the intelligence of communication systems would evolve significantly.

Furthermore, in a classical communication system structure, Shannon’s information theory foundation is based on the *asymmetry of communication* [39]. In other words, data at the receiver cannot be created in an *ex nihilo* fashion. Subsequently, to grant the receiver information generative and reasoning capabilities, one can reframe the communication problem as one in which a pattern or structure (which may or may not be repeated) needs to be realized or constructed in different instances within a limited time duration. This new definition enables reducing the asymmetry between the transmitter and the receiver as we envision for a teacher and apprentice. Here, we ask two important questions that must be answered in order to realize the prior semantic communication definition:

- How should the teacher represent information minimally, without jeopardizing the apprentice’s understanding of the representation?
- What are the steps needed from the apprentice to reason over the received semantic representation and make logical decisions out of it?

Answering these questions necessitates defining the communication problem on the premises of human-like thinking, i.e., transforming the bit/data pipeline information exchange into a knowledge-driven semantic conversation. Answering those questions also requires gradually constructing a semantic language between the teacher and the apprentice. In Fig. 7, we showcase how a teacher and apprentice perform three different transmissions over the course of converging towards a mature semantic language. Fig. 7 also shows the mechanism used by the teacher to *explain* the meaning of conveyed representations. In Fig. 7, we assume that the teacher has the reasoning capabilities needed to *extract* the semantic content elements, and map them to a proper semantic representation. The details of acquiring this skill are discussed in Section IV-A. During the first transmission, we can see that the teacher complements the semantic representation of choice with raw messages. This enables the apprentice to recognize the variability of the transmitted semantic representations with respect to the apprentice’s relevant knowledge. Here, such combination of raw messages and representations that are used during the pedagogic process are dubbed, *semantic didactics*. Then, we can see that, after the apprentice has leveraged causality via queries<sup>3</sup>, semantic

<sup>3</sup>Queries are interrogations posed by the apprentice to learn more information about the causal and associational structure of a representation. In essence, this mimics the classroom learning process, whereby a student asks questions to build their knowledge base on a subject. Queries can be interventions or counterfactuals, and will be detailed in Section V.

didactics contain less raw messages during transmission 2. That is, transmission 2 mostly relies on previously acquired semantic representations, while marginally complementing the unseen representations with raw data. Finally, in transmission 3, we can see that the apprentice has posed the majority of their queries so far, and their relevant contextual knowledge enables them to understand the representations used by the teacher. In this case, we can see that in transmission 3, the teacher solely relied on semantic representations to transmit the information. Consequently, the apprentice can generate the intended image by reasoning over the semantic representations conveyed.

Furthermore, we can observe that successfully representing a peculiar structure by the teacher while successfully understanding it at the apprentice depends on a number of factors: a) the relevant contextual knowledge of the apprentice vis-à-vis the current message to be conveyed, b) the capability of the apprentice to represent the current message while leveraging the relevant contextual knowledge, and c) the level of synchronization between the teacher and apprentice, i.e., how well they are acquainted with their representations. For instance, in the example of Fig. 7, the apprentice has seen pictures of dogs before from other teachers (from previous information exchanges). Nonetheless, the apprentice does not know how to represent Max, and thus two scenarios are plausible here: 1) The apprentice needs to attempt to represent Max given that they have represented a dog before, however this depends on the richness of their knowledge base as well as their reasoning capability, 2) The teacher represents Max for the first time while also complementing the representation with classical data (via semantic didactics); after several transmissions, the apprentice learns how to construct the realization of Max minimally. From our example, we can see that the apprentice asked the teacher via query for more information, proving that scenario 1 might have led to further errors (if such information was not properly requested and delivered). Then, in transmissions 2 and 3, we observe that, within the overall message to be conveyed, Max’s structure is *well-defined* for the apprentice. However, we can also see that the apprentice lacked the representation structure to realize the local park. Here, the apprentice needs to conduct a series of queries to be capable of representing the local park appropriately. Finally, in their last interaction in Fig. 7, we can see that the apprentice has the appropriate relevant knowledge that enables him to interpret the message. Consequently, in such a scenario, minimal representational information is only needed to construct these three well-defined patterns (dog, sunflower, and forehead), already acquired by the apprentice a priori. It is also important to note the following:

- As the language gains more maturity, the communication link’s symmetry increases. As such, in a symmetrical setting, the maturity of the language as well as the enhancement in the reasoning capabilities enable the link to be minimally reliant on the channel and communication resources. That is, as explained in Section II-B3, the apprentice can leverage their knowledge base to under-

stand the cause of particular errors caused by the channel. Also, they can rely on a few representations to generate the remaining content via their computational generative capabilities.

- Our example in this exposition was limited to a visual dog example, nonetheless, this is only for simplicity and for illustrative purposes. The same analogy can be extended and generalized to any *structure* in the data.

Building on our intuition from this example, next, we investigate the necessary measures to migrate from information theory to semantic theory.

### B. From Information Theory to Semantic Information Theory

Fundamentally, as shown in our illustrative example in Fig. 7, while the semantic didactics (a combination of semantic representations and classical data transmitted) are highly dependent on the capabilities of the teacher and apprentice, each unique semantic representation hinges on the *complexity of the task to be described*. A task here is defined as the process of extracting the semantic content elements from a specific datastream, and then mapping each semantic content element to a corresponding semantic representation. In other words, there exists a relationship between the complexity of the task, the number of semantic representations that must be exchanged, and consequently the number of interventions required for an overall comprehension of the apprentice. Nonetheless, given that information theory is solely focused on characterizing uncertainty in the message rather than meaning, it is not the right tool for capturing or characterizing all of these semantic-based concepts. It is thus necessary to propose a suite of novel equivalent fundamentals that build on information theory, and that are tailored to reasoning-driven semantic networks. In essence, to characterize the complexity of a semantic language, we first investigate the way Shannon defined the concepts of information and entropy. This is fundamentally important as designing any communication network and evaluating its performance depend on the way Shannon defined these concepts. Building on these concepts, we similarly draw parallels from classical information theory to future semantic information theory. Such equivalent preliminaries and fundamentals enable us to understand the operation, design, and performance of future semantic communication system in terms of building a semantic language and reasoning over semantic representations.

#### 1) Information: From Uncertainty to Semantic Substance:

Shannon defined information based on its combinatorial nature rather than on its meaning. In other words, under information theory, information characterizes what a transmitter, based on the way bits flow from one medium to another, *could* say, bearing in mind source, channel, and destination errors. Hence, Shannon defined the amount of information according to the number of combinatorial choices that can be considered by the transmitter, i.e., in the simplest cases, the logarithm of the number of available choices in which one can transmit

the message. If we let  $X = x$  be the observed event by the transmitter, the information of observing  $x$  is defined as [40]:

$$I(x) = -\log_2[p(x)], \quad (1)$$

The logarithm function here enables measuring information in an additive fashion with respect to the number of states  $p(x)$  of the system. For instance, if the transmitter wishes to transmit the number 5, 555, 555, roughly 23 bits are needed to represent this rather simple identical sequence of numbers. Interestingly, one can notice that, if a person were to say the exact same sequence to another, they would leverage the existing pattern and say: *the number 5, 7 times*. Leveraging this pattern would inherently minimize the number of bits needed from 23 to almost 10 bits. Alternatively, if a string of numbers with the same length, such as 540, 938, 1326 were to be transmitted then, again, the transmitter needs 23 bits to send them across the channel. In contrast to the previous example, given that this sequence of bits does not exhibit any pattern, as humans we ought to memorize it in our short-term memory before we convey it to the listener<sup>4</sup>. Notably, from these simple, but insightful examples, we can make the following key observations:

- When the message to be transmitted has inherent patterns that can be leveraged, relying on a reasoning mechanism enables to transform the process from a mere recovery process, to a human-like exchange that can potentially minimize communication resources substantially.
- While it was fairly easy to leverage the identical series of numbers in the first example, the second example seemed like an arbitrary sequence of numbers that does not exhibit any regularity. However, this was because the examples were taken *out of context*. The sequence could actually map to a phone number whereby the first three digits map to the state of Virginia area code, the following three digits map to the central office, and the last four digits map to the line number. Here, attempting to learn this task might not enable us to directly minimize communication resources. That is, in analogy to our dog example in Fig. 7, the apprentice would require the teacher to complement the representation with raw data so that the context and the representations are acquired. Then, moving forward from that point, the apprentice is cognizant of context, and, thus, they can utilize their knowledge base to correct any *structurally inconsistent errors* in the semantic representations that map to a phone number or to information that should *logically* accompany such numbers. This will also enable the communication link to be more minimal by not communicating every digit.

If we further apply the insights we have made on the illustrative example of Fig. 7 we can see that, in a classical communication

<sup>4</sup>It is important to note that while the structure in the exhibited example is fairly simple and can be captured via source coding to minimize redundancy, it was used here for exposition purposes in conveying the existence of *long-term and recurrent* structure in the datastream.



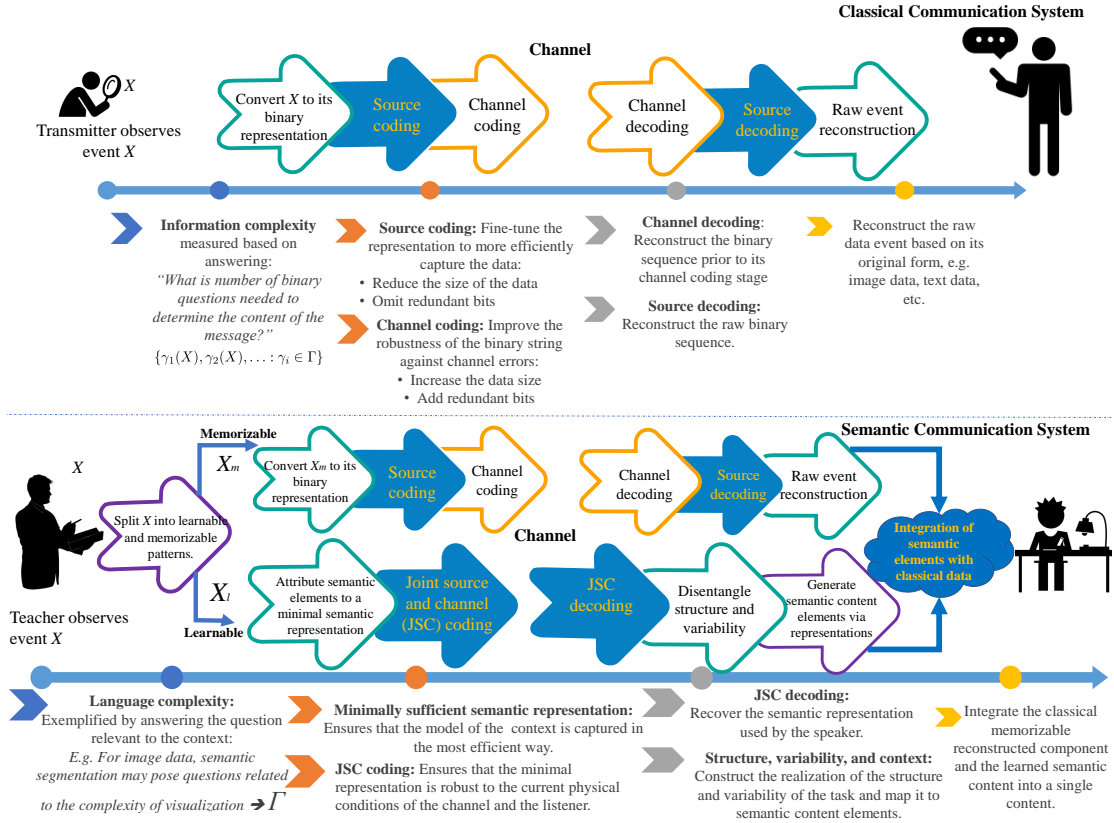


Fig. 8: Illustrative example showcasing the transformation of an E2E communication system from a traditional setting to a semantic-aware one.

system, a  $3,300 \times 4,800$  image will still carry identical information regardless of the simplicity, redundancy, or complexity of the semantic content elements in the picture. This showcases how “information” in Shannon’s definition characterizes what the pixels *could* represent in binary. Meanwhile, if one focuses on the semantic content elements of the image, the information transfer becomes a function of the complexity of the message as well as the maturity of the language established. That is, in the case of a message with an *obvious* structure, like the repeated series of numbers, as well as a mature semantic language, the information transfer process becomes minimal, and it can be easily communicated. Meanwhile, in the case of a complex message, and a weakly-established language between the teacher and the apprentice, the information transfer process may *first* waste communication resources to establish the language. Then, the teacher/apprentice benefit from the knowledge base and language built to ultimately reach an efficient, generalize, and minimal link.

Consequently, Shannon’s “information” is technically quantifying what has been historically known as *syntactic information* [41]. Syntactic information quantifies how much the knowledge of the state of one system reduces the statistical uncertainty about the state of the other system, possibly at a different point in time. In a communication setting, Shannon’s information theory measures how much knowledge about the transmitted datastream reduces the statistical uncertainty about

the state of the received datastreams. While these statistical correlations between the transmitted and received datastream are important, Shannon’s notion of information does not consider what such correlations mean [41]. In contrast, in order to introduce a meaning and context to the definition of information one must: a) Highly correlate the link between the definition of information and the overall goal of the system (if and when such unified goal exists), and b) Capture the relationship between information and causality. Subsequently, we first define the concept of “*semantic information*” for goal-oriented semantic communication systems and then generalize this concept for all semantic communication systems:

**Definition 1.** A particular datastream  $x$  is said to be rich in semantic information if transmitting it improves the system’s ability of pursuing its specified goals.

For example, the goal of a system of digital twins [42] is to guarantee a high synchronization between the physical and cyber twin. Thus, the information exchanged is deemed to be a valuable semantic if it is capable of enhancing the performance of the digital twin in terms of reliability and synchronization. As discussed, in many use cases, a unified goal might not exist. In such settings, we extend the definition of semantic information to what follows:

**Definition 2.** A particular datastream  $x$  is rich in semantic

information if it improves the reasoning capabilities of the teacher and the apprentice to ultimately expand the semantic language built between them.

The concept of semantic language is explained, in depth, in Section IV-B. Essentially, Definition 2 captures the fact that a datastream is void of information if it does not contribute to enhancing the reasoning capabilities of the teacher and the apprentice. The improvement in reasoning capabilities can be achieved by leveraging causality, and properly identifying the root-cause of new semantic content elements, based on previously observed ones. Also, similarly to humans, when our learning improves, our capability to express our understanding also improves. Hence, a semantic language is a key element that captures the capability of a radio node to ultimately understand the meaning of information. As such, the previously established definitions of semantic information are the initial guide in the process of defining a semantic language and its complexity, compared to entropy. Next, we investigate the information theory fundamentals of entropy, then, we further discuss the necessary measures to define its counterpart in semantic communication systems.

2) *From Entropy to Language Complexity*: The overarching role of entropy is to characterize the uncertainty of a micro-state on the macro-state of an overall considered physical or statistical system. Stated differently, statistical entropy is proportional to the number of “yes/no” questions that must be asked, to determine the micro-state, given that the macro-state is known. In classical communications, the micro-state or atomic unit is a “bit”, consequently, if we let  $\Gamma$  be the set of all possible binary functions on a domain  $X$ , characterizing the entropy requires asking the following question: *What is the optimal sequence of yes/no questions needed to construct*:

$$\{\gamma_1(X), \gamma_2(X), \dots : \gamma_i \in \Gamma\}, \quad (2)$$

with the goal of perfectly recovering  $X$  from the shortest sequence of binary answers [43].

Furthermore, based on the previous definition of *information* in (1), answering this question requires computing the entropy over the domain  $X$  in the classical sense, as follows:

$$H(X) = - \sum_{i=1}^n p(x_i) \log p(x_i). \quad (3)$$

From (3), we observe that the entropy varies with respect to the freedom of choice of the transmitter. That is,  $H$  tends to 1 when a single probability measure  $p(x_i)$  is equal to 1, and all other probabilities are equal to 0. In this case, it means that one can *certainly* represent the message deterministically. However, the entropy is independent of the message itself. We can also observe that, when we are in an equiprobable setting, then entropy is at its highest. Meanwhile, if the states are not equiprobable, compression becomes more possible [41], and thus fewer yes/no questions need to be asked to identify a particular input of datastream. Nonetheless, in all cases, i.e. whether compression is considered or not, entropy remains a

function of the choices of the transmitter and not a function of the message itself.

Interestingly, if we consider the way humans communicate: when trying to explain something, we do not attempt to do so while considering the number of yes/no questions required to explain a particular story. Such yes/no questions fail to characterize the *structure* of the story itself. Meanwhile, we tend to focus on the most *meaningful* information that constitutes the *core* of the story. To do so, the set  $\Gamma$  needs to be restructured so that it consists of questions regarding the content complexity (which highly correlates with the semantic content elements) rather than uncertainty (which is devoid of it). One can for example let  $\Gamma$  represent the visual semantic entropy [44], [45], whereby the set of queries would indicate the presence or absence of an entity, and its relation to the other represented entities. From our Fig. 7 example, a simple question to semantically represent transmission 1 would be “Is there a dog at the center of the picture?”. Meanwhile, semantically representing transmission 3 requires asking the question “Is there a dog with a sunflower on its forehead?”. We can see that the latter has a substantially higher reasoning complexity as multiple entities are involved, and they are related to one another via a certain function, and there exists a function between them (the sunflower is exactly on the forehead). Indeed, these questions enable us to characterize the complexity of a task to be semantically represented, nonetheless, the problem is that the set  $\{\gamma\}_{\gamma \in \Gamma}$  cannot be computed tractably in this case.

In essence, we need an alternative to entropy that enables exchanging questions between the transmitter and receiver to ultimately capture structure in the information and reflect the semantic content of the datastream. Consequently, this exchange reflects the necessity for engineering a semantic language between the transmitter. Here, *the concept of language complexity in a semantic-based system is the equivalent of entropy in a classical one*. Next, based on the nuts and bolts we have elucidated to extend classical information theory to semantic information theory, we will investigate the fundamentals of semantic representations and languages for semantic communication systems.

#### IV. BUILDING SEMANTIC COMMUNICATION SYSTEMS: REPRESENTATIONS AND LANGUAGES

As shown in Fig. 8, classical communications starts by first processing the raw data and converting it into binary representation. Then, this binary representation is source coded to minimize the number of redundant information bits. Then, this resulting datastream is channel coded to improve its robustness against the adverse channel conditions. At the receiving end, such processes are mirrored to ultimately reconstruct the raw data message originally transmitted. As discussed in Section II, all of these processes use Shannon’s concept of “information”, viewed as “uncertainty”, and, thus they do not capture any core structure in the data. Moreover, they are not able to attribute any meaning to the datastreams processed. In this section, we investigate the concept of structure and variability in the data

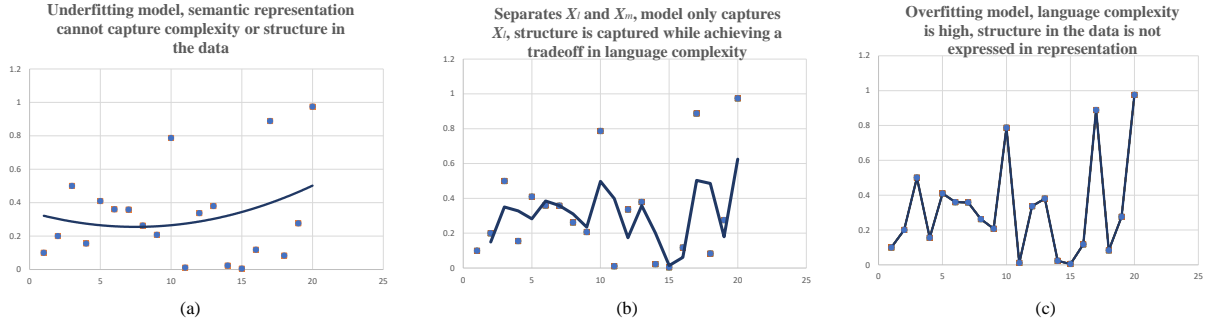


Fig. 9: Illustrative example showcasing the difference between (a) an overly simple model, underfitting over the raw data, (b) a reasonable model that separates  $X_l$  from  $X_m$ , captures the semantic data, and leaves out random data, and (c) an overfitting model that captures all the characteristics of the data without exhibiting reasoning capabilities, and without capturing structure.

and, then, we highlight the necessity of splitting the data into *learnable and memorizable* patterns. Then, we will discuss the properties of a semantic representation and language. We will further gradually explain the semantic communication system building blocks shown in Fig. 8 while we explore these novel concepts in this section.

#### A. Language pre-processing

In any communication system, the first step is to convert the highly dimensional raw data into binary data suitable for communication. While in classical communications this is customary and followed by source and channel coding, in a semantic communication network, we must scrutinize the data, and attribute a representation to every major structural part in it. Nonetheless, if one considers the whole datastream as a bulk, the learning process, i.e., identifying every semantic content element and crafting its corresponding semantic representation, becomes highly inefficient. This is due to the fact that raw data contains a lot of purely random information. This is problematic because such random information increases the complexity of the built language (which we will soon formally define), yet contributes to a language with spurious semantic representations. That is, one might think that the semantic language complexity stems from an inherently complex semantic structure (e.g. a complex high-dimensional hologram that must be described via a semantic language). In reality, the deceiving complexity measure here stems from an abundance of random information in the datastream.

Hence, for the sake of efficiency, prior to learning a semantic language, we must first separate two parts of the data, learnable and memorizable, as defined next:

- 1) The *learnable part  $X_l$  of the data*: is the core of semantic information. That is, performing *reasoning* on such learnable data points allows the radio node to capture the inherent structure in each semantic content element. Thus, reasoning over  $X_l$  eventually leads to a meaningful semantic language with non-spurious models. One can use  $X_l$  as an input to build a proper semantic language. This is shown in Fig. 10.

- 2) The *memorizable part  $X_m$  of the data*: is the one that is governed only by pure randomness. For example, these can be details in an image that do not contribute to any semantic content element. As such, attributing a semantic representation to  $X_m$  is a highly *complex* learning process, but a very *simple* memorization process. In other words, learning the structure of a purely random and structure-less datastream is *complex* process and does not yield meaningful semantic representation. Thus, allocating computing resources to *learning* random information is a wasteful process, thus  $X_m$  must be transmitted using classical communication resources.

Furthermore, as shown in Fig. 9, separating data into memorizable and learnable components is crucial. If one attempts to learn all the data points as shown in Fig. 9 (c), the resulting semantic representations learned will be spurious. While the curve here “fits” the data, it does not learn any inherent structure.<sup>5</sup> In this case, reasoning on all the datastream, will not reflect the true semantic content elements of the data. This will lead to an inherently poor semantic language that cannot express structure in the data. Meanwhile, if one only captures the learnable part of the data as shown in Fig. 9 (b) to build a language, the learning process would not overfit. It will capture all data points as its goal is to capture structure and not *memorize* information. It is also important not to have an overly simplified language (Fig. 9 (a)) that will fail to yield a minimal representation that characterizes the structure of data. We will further revisit Fig. 9 when investigating the semantic language complexity. It is important to note that, in classical ML frameworks, the data is not separated or pre-processed as the goal of ML is fundamentally different from the overarching goal of semantic communications. In essence, in ML the goal is *learn* the data and then leverage this learned-model to make some predictions. Nonetheless, in semantic communication the goal is to perform *reasoning* and dissect the

<sup>5</sup>In fact, the overfitting scenario shown in (c) is what source coding and data compression attempts to do, whereby they just aim to summarize data and then reconstruct it bit-by-bit. The compression does not identify any structure

structure of semantic content elements. A semantic language and a knowledge base cannot be built via a series of models and bulks of data (see Fig. 1), but they can only be built via organized knowledge (e.g. causal reasoning schemes and languages).

Moreover, as shown in Fig. 8, in contrast to a classical communication system, the first step in semantic communication systems is to split the observed data into learnable and memorizable parts. Then, the path followed by  $\mathbf{X}_m$  is identical to the one of classical communications. To understand this more clearly,  $\mathbf{X}_m$  mimics data points in our daily lives that are *wasteful* to learn. For instance, when trying to call someone, the digits following the area code are purely random, they do not follow a specific pattern nor do they result from a causal event. Thus, in such scenarios, our mind operates like a bit-pipeline that captures the information without performing any reasoning over it. Note that the separation between  $\mathbf{X}_m$  and  $\mathbf{X}_l$  is not unique, which makes this process of separation challenging. Here, one can resort to ML schemes to perform this categorization [46], [47], or can leverage causal models as discussed in Section IV. That said, these disentanglement methodologies remain outside of the scope of this work.

We can also see in Fig. 8, that for the semantic path,  $\mathbf{X}_l$  first undergoes a *reasoning* process whereby its semantic content element and representation are identified. Then, every semantic representation, and ultimately the semantic language, undergoes joint source and channel coding. This is intrinsically different to classical communications whereby the *source data* undergoes joint source and channel coding. Here, the goal of joint is to leverage the correlation between multiple semantic content elements (or even for multiple transmitters) for cooperative transmission [48]. In fact, as laid out in [49] and [50] joint source and channel coding enable executing specific actions (in a goal-oriented scheme) directly rather than recovering the entire source messages. Similarly, in a broader semantic communication system, joint source and channel coding schemes enable re-generating (at the destination) the specific semantic content elements rather than the whole datastream. This is necessary as the receiver no longer decodes one bit stream at a time, but rather generates one semantic content element (meaning) at a time.

### B. On the Structure and Complexity of a Semantic Language

Thus far, we have pre-processed the data and separated  $\mathbf{X}_m$  from  $\mathbf{X}_l$ . In the sequel, and since  $\mathbf{X}_m$  will undergo the path of classical communication, we focus on building a language from input  $\mathbf{X}_l$ . The concept of a semantic language can be seen in Fig. 10. We first define a semantic language and then scrutinize its structure and variability to ultimately assess its complexity.

**Definition 3.** A semantic language  $\mathcal{L} = (X_{l,i}, Z_i)$ , is a dictionary (from a data structure perspective) that maps the learnable data points  $X_{l,i}$  to their corresponding semantic representation  $Z_i$ , based on the identified semantic content elements  $Y_i$ .

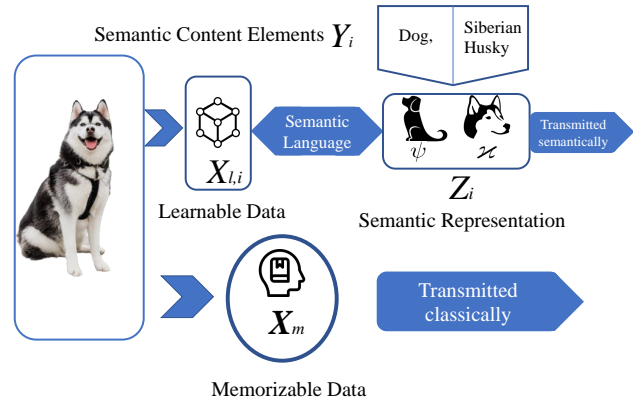


Fig. 10: Illustrative example showcasing the concept of semantic language.

Representation  $Z_i$  must be efficient in inducing the apprentice to *generate* the originally conveyed task or semantic content  $Y_i$ . Thus, language  $\mathcal{L}$  re-purposes the apprentice’s task from a mere *reconstruction* mechanism to a *generative and automatic* process. That is, the language re-purposes information transfer via the convergence of computing and communications in semantic systems, thus yielding a *generative apprentice and a dictating teacher* as follows:

- **Generative Apprentice:** The apprentice relies on their *computing resources* and reasoning faculties to *generate* content from a *representation*. This mimics our imaginative experience when someone mentions a particular term. That is, when someone says “flower” we can “imagine” and “recreate” what a flower looks and feels like, based on its representation. For instance, in an XR setting, if the transmitter is willing to augment the metaverse with new XR content which is a flower on the sidewalk, they have to modify every bit surrounding the metaverse 3D space so as to reconstruct a flower on the sidewalk. Meanwhile, in a semantic setting, the teacher “describes” to the apprentice that they need “add a flower, on the left corner of the sidewalk”. Based on their previous language, the apprentice maps findings in their knowledge base to what is “flower”, “left”, and “sidewalk”, to generate it via its computing resources.
- **Dictating Teacher:** The teacher does not have to exhaust multiple communication resources to convey a message in a *bit-by-bit* fashion. Instead, the teacher must: a) identify the semantic content elements in the data, i.e., learn the meaning in the datastream (via their computing resources), b) based on this meaning, the teacher must attribute an existing representation or develop a new representation, and add it to their language, and c) convey the message based on the representations and language developed.

Essentially, each semantic representation within language  $\mathcal{L}$  must be able to:

- 1) Characterize the *structure*,  $\psi$ , that describes the data collectively. For example, this could be the main semantic

element of an image. For instance, if a semantic representation describes a German Shepherd, the common denominator between the German Shepherd and all dog breeds represents the *structure*. Naturally, this concept applies to any type of data, and not just images.

- 2) Characterize the *variability*,  $\varkappa$ , of the individual data points with respect to the shared structure. For instance, taking the example of a German Shepherd, variability characterizes distinctive features of a German Shepherd compared to other dogs.

Furthermore, we can define the *complexity of the language*  $\Gamma(\mathcal{L})$ . The goal of this complexity is to characterize the difficulty of identifying and learning the semantic content elements within  $\mathbf{X}_l$ . This process ultimately stems from the inherent structure of the data (which we will formally define shortly). This complexity also captures the difficulty of developing a language for this semantic content, i.e., a semantic representation for every semantic content element. Our framework has been built by borrowing the definition of dataset complexity in transfer learning in [51] and [52]. In the following proposition, we characterize this language complexity.

**Proposition 1.** *The complexity of a specific language  $\mathcal{L}$  adopted among a teacher and apprentice pair is given by:*

$$\Gamma(\mathcal{L}) = \min_{p(\mathbf{Z}|\mathbf{X}_l)} L_{\mathcal{L}}(p) + K(p). \quad (4)$$

Here,  $L_{\mathcal{L}}(p) = \sum_{i=1}^N -\log p(Z_i|X_{l,i})$  is the cross-entropy loss, and  $K(p)$  is the Kologomorov complexity of the distribution  $p(\mathbf{Z}|\mathbf{X}_l)$ .

From Proposition 1, we observe the following:

- 1) When learning random semantic-agnostic representations for the data, the complexity  $\Gamma(\mathcal{L})$  will become very high. In essence, this is a result of the fact that the majority of the data *lacks structure*. Thus, in that case, the “task” of concern should be memorization rather than learning (using a language to communicate random information is a highly inefficient task). Hence, the teacher must prune data points from  $\mathbf{X}_l$  and attribute them to  $\mathbf{X}_m$ , i.e., the pre-processing phase must be revisited. This will in turn minimize the complexity  $\Gamma(\mathcal{L})$  and makes learning the language more efficient.
- 2) The language complexity is a metric that replaces the classical entropy. It is a function of the cross-entropy loss, whereby the fitness of the model is captured as well as the Kologomorov complexity of the data. Here, in contrast to entropy which only characterizes the uncertainty and the freedom of choice at the transmitter (which we highlighted in Section III-B2), language complexity measures the complexity of the raw data and the language model at stake when communicating a particular message.
- 3) Unlike Shannon’s information-theoretic perspective whereby the encoding for a message is predetermined by the randomness of the source transmitting it, Kolmogorov’s complexity enables us to characterize

the *individuality* of the semantic content elements in the message to be conveyed. In fact, Kolomgorov’s complexity  $K(x)$  is a measure of the shortest effective binary description of  $X$ . In other words, it characterizes the methodology that enable the apprentice to *autonomously generate* the semantic content elements based on the conveyed semantic representations.

Essentially, the probability distribution function (PDF)  $p$  in Proposition 1 constitutes the learned model of the semantic representation of a particular semantic content element, based on input  $\mathbf{X}_l$ . Consequently, it is necessary to characterize the tradeoff between the ML loss achievable by this model for a language  $\mathcal{L}$  and its complexity. This tradeoff essentially is a *measure* that indicates to the teacher/apprentice: a) the maturity of the established language, based on how expressive it is with respect to semantic content elements, and b) the language complexity which will indicate the level of reasoning required, as well as the amount of computing resources needed to achieve this language. This tradeoff is captured via the *structure function*:

**Definition 4.** *The structure function achievable by a model  $p$  for a language  $\mathcal{L}$  is given by [51]:*

$$\Psi_{\mathcal{L}}(t) = \min_{K(p) \leq t} L_{\mathcal{L}}(p). \quad (5)$$

The structure function  $\Psi_{\mathcal{L}}(t)$  will be zero for sufficiently high complexity. In particular, after all the shared structure is captured and characterized, the only methodology that minimizes the loss in it is by memorizing the variability of the leftover data bits, leading to an overall high structure function.

To solve the optimization problem in (5), one can rewrite (5) with respect to its associated Langragian:  $L_{\mathcal{L}}(p) + \lambda K(p)$ , where  $\lambda$  is the Langrangian multiplier. Furthermore, taking the minimum over  $p$  leads us to obtain a family of complexity measures for our language  $\mathcal{L}$  parameterized by  $\lambda$ :

$$\Gamma_{\lambda}(\mathcal{L}) = \min_{p(\mathbf{Z}|\mathbf{X}_l)} L_{\mathcal{L}}(p) + \lambda K(p). \quad (6)$$

(6) is nothing but the Legendre transform of the structure function  $\Psi(t)$  as a function of  $\lambda$ . Solving (6) can be performed by increasing the complexity  $K(p)$  of the model until the return obtained is smaller than the constant  $\lambda$  chosen by the teacher.

Moreover, the model  $p(\mathbf{Z}|\mathbf{X}_l)$  is considered a Kolmogorov sufficient statistic of the language  $\mathcal{L}$  if it minimizes the complexity formulated in Proposition 1. This reconfirms our initial rationale in Fig. 9 motivating the need for a model that acts as the smallest statistic that can characterize a particular semantic representation. Thus, to ultimately reach a semantic representation that is minimal and efficient, the statistic  $p$  must be able to learn and build the language  $\mathcal{L}$  without squandering computing and communication resources to model random, non-semantic information. Instead, the minimally sufficient statistic  $p$  should be able to characterize the valuable and semantic information within the data, i.e.,  $\mathbf{X}_l$  as previously shown in Fig. 9.

### C. Why Semantic Languages? Why not Natural Languages?

In Section. IV-B, we have discussed the necessity to categorize the *learnable and memorizable* data patterns in order to create a more efficient E2E semantic communication system. Then, based on the learnable data categorized, the teacher and apprentice need to create a language  $\mathcal{L}$  that yields semantic didactics that have a particular structure and variability. Essentially, the complexity of learning  $X_l$  mainly depends on the structure function of choice. A semantic language with an extremely high complexity  $\Gamma(\mathcal{L})$  can be a result of:

- 1) A poor separation between  $X_l$  and  $X_m$ . This leads to an  $X_l$  with an inherently poor and unlearnable structure function. Thus, in this case, the teacher needs to revisit the methodology they are adopting to categorize  $X_l$  and  $X_m$ .
- 2) Given a particular  $X_l$  that possesses a rich structure function, the teacher needs to revisit the methodology adopted to solve the optimization problem in (6). In this case, the high complexity can be a result of not finding a *proper sufficient statistic of the language  $\mathcal{L}$* .

We have so far investigated the fundamental notions of a semantic language and its complexity. We have discussed how a language re-purposes the information transfer with a *dictative teacher* and a *generative apprentice*. Here, given that this information transfer mimics human conversation, one might think that natural languages are the solution that can link the teacher and the apprentice. This misconception has been pervasive in the semantic communication literature [26], [53], and [54]. Nonetheless, equipping radio nodes with a natural language constrains the information transfer process with various challenges that include wording and deterministic syntax rules. Also, a semantic representation must be characterized with features that are fundamentally different than words (which are the atomic unit of a natural language). Next, we formalize our ideas further by highlighting the common denominator and contrasting features between a natural language and a semantic-centric language. Then, we highlight the key characteristics of a semantic language and representation.

1) *Properties of a Natural Language*: A semantic representation is the atomic unit of a language. Thus, to formalize the concept of a semantic representation, we first contrast the differences between a natural language and a semantic language. In general, a *natural language* must have the following properties:

- **Syntax**: This is a system of rules constructing the possible grammatical and acceptable sentences out of words (symbols), and determining their sequential arrangement to create a well-formed sentence (expressions).
- **Semantics**: This is a system that attributes a meaning to a well-formed sentence (expression) constructed according to a particular syntax.
- **Pragmatics**: This is a system that specifies how the semantically constructed syntax in a language can be used. In other words, pragmatics are the foundation of the

context-dependent features of a language. For example, if someone asks the question, “Are you wearing your seatbelt?”, this is a sentence that urges the passenger to exercise cautiousness, although the word “cautiousness” has not been used in the sentence.

Indeed, every natural language requires the aforementioned properties to enable a smooth communication between human beings. Nonetheless, while semantic communication systems should mimic human communication, their goals and operation will be slightly different. Next, we emphasize the distinctive features of a semantic language, and we highlight why it is fundamentally different from a natural language.

2) *Goals and Properties of a Semantic Language*: Building a language between the teacher and the apprentice in a semantic communication system does not necessarily require it to be a natural language. Thus, in what follows we elaborate the standing of a semantic language with respect to the key properties of a natural one:

- **On syntax**: Relying on syntax and a set of grammatical rules to construct meaningful expression of semantic representations defies the initial purpose of semantic communication systems. In other words, the *structure* of syntax is characterized by the properties of the deterministic rules set. Meanwhile, the goal of a semantic language is to characterize the causal and statistical properties of the datastream via a semantic representation. Thus, the core of a semantic representation is the structure of the data to be transmitted and not a set of deterministic rules. Hence, adding syntax will only add an overhead of deterministic rules to the established semantic language.
- **On semantics**: Inherently, as the name suggests it, the goal of a semantic language is to transmit information by focusing on its meaning. Hence, as will be elaborated next, a semantic language aims to be simpler, by avoiding the strict formalities of a natural language. This is similar to when two people are used to each other’s vocabulary and avoid formal language. This ultimately *minimizes back-and-forth messages exchanged*.
- **On pragmatics**: Relying on pragmatics asymptotically requires reading between the lines. While this is a consequence of any natural language, it is a sufficient condition but not a necessary one in a semantic language. In semantic communication networks, pragmatics are highly correlated to the context or general theme of communication related to the current semantic messages communicated. Indeed, this skill improves the “dynamic” reasoning capability of the teacher-apprentice pair. Dynamic reasoning is characterized by generalizing the reasoning performed on one scenario to various different settings (e.g. extract location-based intelligence and use it for tasks such as environmental sensing). However, this skill can only exist when the language has matured. In other words, this skill emerges when the language’s complexity has been minimized while capturing structure, i.e., the problem in 5 has been solved. This also needs to be accompanied

with a known (or non highly varying) context of communication. The sentence “Are you wearing a seat-belt?” urged cautiousness in the language when the “context” of communication is known to be in a vehicle. In fact, the concept of pragmatics justifies why a *consistency in the context of subject* improves the dynamic reasoning capability of semantic radio nodes.

Consequently, given the fundamental differences between natural languages and semantic languages; next, we delineate the key properties of a semantic representation and its respective semantic language (representations are the atomic constituents of a language).

- **Minimalism:** Based on our observations from Proposition 1, a semantic representation needs to be *minimally sufficient*. That is, if a semantic representation is very simple, it has traded-off part of the structure of the message to minimize the communication resources. Meanwhile, a semantic representation that is overly complex, will have an overfitting AI model, i.e., high complexity and low structure. Such a representation would be describing unnecessary random information in the data.
- **Efficiency:** While a natural language does not mandate efficiency, a conversation between a teacher and an apprentice in a semantic communication network must be more efficient than a classical communication paradigm. Efficiency can be measured via the semantic impact (which is a new metric that we propose in Section VI). In essence, efficiency is a measure of the time and communication resources that can be gained from adopting a semantic communication system, in comparison to a classical one. For example, assuming that the goal is to send large XR content to the apprentice (contributing to better quality-of-experience (QoE)), if learning the structure of the data does not ultimately lead to a faster and more efficient generation of XR content at the apprentice; classical communication is a better option.
- **Generalizable:** A specific semantic representation needs to be distribution, domain, and context invariant [55], [56]. This invariance evokes generalizing to *new and unseen, out of distribution, domain, and context* data points. In other words, the teacher must be able to extract particular features within the data whose structure is invariant to the context or domain. Context generalization can be seen in our example in Fig. 7, here, the teacher needs to be able to semantically represent Max regardless of the context Max occurs in, and the correlations that such context might have with the content element. Meanwhile, domain generalization allows the reasoning node to learn a representation and generalize it to unseen target domains. For instance, assuming a semantic representation has been learned from an Internet of Things (IoT) data source within image data, the reasoning node must be able to use this representation to describe the same semantic content element if it appears in the metaverse.

Thus far, we have elucidated the foundational characteristics of a natural language and its contrasting features to a semantic-centric language. In addition to natural languages, there exists various forms of representations that are currently being considered to serve as a candidate semantic representation. In what follows, we overview the potential caveats of such representations.

#### D. On Existing Forms of Semantic Representation

Many efforts in the research community attempted to devise methods or mechanisms that yield a semantic representation of data. We have comprehensively detailed their contrasting characteristics in Table II, and we explain them next.

1) *Natural language processing (NLP):* The goal of NLP is to transcend the capability of a computer and grant them the capacity to understand and speak human language. Recently, a number of works in [26], [53], and [54] have adopted NLPs in designing semantic communication systems. However, such works limited their communication to *text data*, this is why a natural language was sufficient for *describing* the data. In essence, describing complex messages like a hologram, or a high-precision manufacturing command using natural languages leads to more repetitive and redundant communication resources to describe the message. That is, our natural languages fail to “describe” highly complex processes. For example if one tries to use the English language to describe the hologram of a particular entity, it would require a lot of time and sentences, that might not be able to ultimately generate the hologram in the intended fidelity. Meanwhile, a coded language has the capability to do so. Here, the semantic representation would be describing to the apprentice a mechanism similar to that code, which would automate the generation hologram via computing resources on the long-term. As outlined in the previous section, there are fundamental differences between a natural language and a semantic language. NLPs are limited to *wording* in the same way that classical communication are limited to bit-pipeline. Additionally, NLPs require syntax and a set of deterministic rules, which are intrinsically not concerned with the meaning of the data. Hence, while NLPs are a good AI tool for the *service intelligence* of some tools like an automated chat bot or robot, if applied as is, they will fail at the helm of reasoning for low-layer data.

2) *On ANNs:* ANNs mathematically exploit the universal function approximation theorem [57] to improve their generalization capability. Recent works in [58]–[63] considered ANNs as a tool to design a semantic encoder and decoder. While ANNs are a powerful tool for data analytics, they are limited to the statistical structure of the data. Indeed, ANNs are unable to reason the cause, context, or effect surrounding the event or source of the data. Thus, the reasoning capability of ANNs is more or less limited by the statistical nature of the data – *they are not knowledge-driven*. In other words, if the data is not *purely statistical*, an ANN will not be able to capture the complexity of the data nor learn a proper representation. That said, ANNs remain powerful building blocks that are needed in our

Table II: Characteristics contrasting potential approaches to reason semantic representations

Key metric	NLP	ANN	Knowledge Graphs	Topos	Causal Representation Learning
<b>Methodology Fundamentals</b>	Read, understand, and decode human languages in a valuable fashion. Subsequently, use a human language to describe the semantic information contained in raw data.	Perform semantic feature extraction, and then leverage a specific neural network structure to model the task complexity and its subsequent semantic features.	Perform causal discovery of semantic features and represent such features and their corresponding relationships via vertices and edges in a knowledge graph.	Translate current data structures to a novel well-defined morphism, that enables extracting unobserved semantic information.	Learn a representation that can partially expose the unknown causal structure of data. This structure unveils the semantic content elements of the data and their relations, thus characterizing the context.
<b>Generalizability</b>	Medium	Medium	Low	High	High
<b>Minimalism</b>	Medium	Low	Medium	Medium	High
<b>Benefits</b>	<ul style="list-style-type: none"> <li>Easily understandable and decodable for design purposes.</li> <li>Suitable for specific data structures that heavily rely on text data.</li> </ul>	<ul style="list-style-type: none"> <li>Easily integrated with existing ML and artificial neural network models of the E2E network models.</li> <li>Has gained maturity and can be easily reparametrized.</li> </ul>	<ul style="list-style-type: none"> <li>Characterizes intertwined parameters in simplified graphs.</li> <li>Provides the apprentice with a simplified basis for reasoning.</li> </ul>	<ul style="list-style-type: none"> <li>Capable of unraveling unobserved contextual patterns</li> <li>Can perform reasoning beyond statistical boundaries.</li> </ul>	<ul style="list-style-type: none"> <li>Can leverage the concept of interventions and counterfactuals to understand the structure of the data beyond associative logic.</li> <li>The causes of semantic content elements implicitly characterize the context of the transmission.</li> </ul>
<b>Disadvantages</b>	Limited by syntax, pragmatics, and wording.	Limited by statistical relationships within the data.	Can only represent simplified causal graphs, and restricted to the expressivity of a graph.	Many concepts within Topos remain intractable and difficult to characterize.	Limited by posing the proper interventions/counterfactuals at the apprentice.
<b>Major Challenge</b>	Transforms the bit-pipeline problem to a word-pipeline problem.	Limited reasoning capabilities with respect to contextual information that are not limited to statistical relationships.	Fails to characterize highly complex tasks despite its causal structure.	Non-scalable for the overall E2E communication system.	How to concatenate a structural causal model within an ANN that characterizes both statistical and causal properties?

E2E semantic communication framework. In essence, ANNs can potentially be used in modeling variables and parameters that describe purely statistical models, such as the channel, joint channel and source encoding processes, etc. While such processes will play a meaningful role in the ultimate semantic representation of the data, ANNs cannot be the main block in charge of yielding the semantic representation.

3) *Knowledge Graphs*: Several recent works [64]–[66] developed semantic communication systems based on knowledge graphs. The approaches in [64]–[66] adopt knowledge graphs as a result of their explainability and interpretability. However, in a realistic scenario, raw data is not structured into separate and well-defined categories and units. Thus, devising a semantic representation is a task that requires disentangling the raw data from memorable information and scrutinize the structure and variability prior to any causal discovery. Moreover, as previously explained, the learnable data  $\mathbf{X}_l$  must exhibit rich semantic structure. However, limiting the representation model to knowledge graphs limits the expressivity of the model, and, consequently, the capability of communicating semantically complex messages.

4) *Topos*: The works in [67] and [68] used the mathematical framework of topos theory for semantic representation.

Toposes originated from a purely mathematical foundation, in particular, homological algebra and algebraic topology [69]. In essence, this technique translates every data structure by a family of objects in a well-defined topos [70]. Here, a semantic representation stems on from the construction of a category set, a class of objects, types and stacks needs to be constructed from the existing raw data. Nonetheless, adopting topos leads to multiple challenges: a) Deploying them requires a major overhaul on all existing machine learning frameworks running at any layer in the network stack. That is, they need to be re-characterized and re-defined in the family of objects in the topos. Thus, topos cannot easily be integrated with co-existing AI frameworks in the network stack, and b) Topos may have high computational complexity and intractability when handling raw data at the teacher’s side prior to assigning it with a particular semantic representation. That said, despite these challenges, topos can play a role as specific building blocks of the E2E network, as they are capable of characterizing a “many-world view” interpretation in certain settings. For example, one can transcend basic probabilistic quantities in information-theoretic settings or semantic-aware ones to a cohomology class (or a cocycle) in toposes [69].

Having concretely defined the tenets of a semantic lan-



guage, and explained the existing frameworks for semantic representations, we can now move our focus to the reasoning aspect, which is one of the most fundamental building blocks of a semantic communication network. Hence, in the next section, we first dwell on the concept of reasoning via causality. Then, we investigate the characteristics of causal representation learning that can lead to a minimal, generalizable, and efficient semantic representation and language.

## V. REASONING VIA CAUSALITY FOR SEMANTIC COMMUNICATION SYSTEMS

### A. Motivation and Preliminaries

As shown in Table II, existing generative methods have their own benefits and challenges in contributing to a proper semantic representation. In particular, one of the most fundamental challenges and drawbacks surrounding most of the current methods is that they excel at finding associations in independently and identically distributed (i.i.d) settings. However, real-world data that is communicated between two wireless nodes is often not i.i.d. For example, this confusion might take place when Max’s representation in Fig. 7 is placed in a datastream that may be negatively correlated with its presence, i.e., background content that may negatively correlate with Max’s data structure [71]. In this case, a) the teacher will experience a difficulty in extracting the semantic content elements and b) the apprentice will face a difficulty in understanding the representation communicated by the teacher. Thus, adopting the previously surveyed AI techniques (see Table II) to extract or generate a semantic representation faces multiple challenges:

- 1) Existing frameworks (surveyed in Section IV-D) cannot adapt to distribution shifts, i.e., the ability to generalize from one data point to the next, when sampled from a particular distribution that exhibits shifts. This is because such frameworks are mostly training and/or bias dependent. Thus, such frameworks exhibit *weak vertical generalizability*.
- 2) The aforementioned techniques do not possess the reasoning capabilities needed to draw logical conclusions out of datastreams, that go beyond generalizing to out of domain/distribution/context. The ability of the teacher/apprentice to perform associative logic and draw logical conclusions is a unique and necessary feature of a semantic communication system, that cannot be achieved by the previously discussed techniques.
- 3) The various frameworks investigated in in Section IV-D cannot adapt to novel domain and context settings. As a result of this caveat, such frameworks exhibit *weak horizontal generalizability*. This is usually exhibited when the teacher and apprentice are communicating in a totally new setting, yet this setting exhibits new representations that have share analogous data patterns to previously exchanged and learned ones.
- 4) Existing AI frameworks generally have a large computational complexity. This drawback is more pronounced for

tools such as knowledge graphs and Bayesian machines. Consequently, adopting such frameworks to build future semantic communication networks may not be scalable.

To remedy these technical shortcomings, it is necessary to move from AI-augmented systems that are information-driven, to knowledge-driven or reasoning-driven systems as shown in Fig. 1. In particular, examining the “*why and how?*” of the semantic content elements in an observed datastream or event, enables distinguishing between the set of meaningful information and the statistical accidents (e.g. identifying “Type I error”). Essentially, a *causal inference* framework consists of two parts: causal identification or discovery and statistical inference. In a semantic communication setting, *causal identification* is the process of identifying the root causes of semantic content elements in the data. Furthermore, *statistical inference* (in a causal-semantic setting) is the process of attributing a minimal semantic representation to that causal structural model. These two processes together create a semantic communication system with causal logic. Moreover, to efficiently perform *causal identification* the following must be considered:

- 1) If a radio node cannot characterize a correlation between the datastream to be conveyed and the established knowledge base (created through an accumulation of previous causal information), then the datastream does not exhibit any causality. Thus, this datastream contains pure random information, and no structure, and, thus, it is more efficiently transmitted classically.
- 2) If a correlation exists, between the datstream and the knowledge base, and there is a unique causal model that can rationalize this correlation, then true causality exists and it can be characterized via a representation.

After evaluating the previously established correlation, the goal of *statistical inference* is to asses the overall correlation and causality, as well as the corresponding stochastic changes in the system stemming from this causality. Furthermore, analyzing the reasoning behind the data pattern requires investigating the three levels of causal hierarchy: a) associations, b) interventions, and c) counterfactuals [72] (which will be explicitly defined later in Section V-B). Notably, the concept of ladder causation was first established in [73] and is also called the Pearl’s Causal Hierarchy. These conceptual levels constitute the foundational pillars of causal reasoning that we propose to use for equipping the teacher and the apprentice with a rationale beyond the statistical association between variables as shown in Fig. 11. Mathematically, these are logical queries that are posed via conditionals (for associational logic) or so-called “do operators” (for causal logic). The “do operator” is a mathematical operator that is performed on a causal model and its corresponding causal network which describe the intertwined cause-and-effect relationships in a particular datastream. It involves two tasks, a) setting a variable to a specific value (e.g. setting the parent node of branch within the causal network, to a certain value), and b) removing a particular branch in the causal network that contribute to a

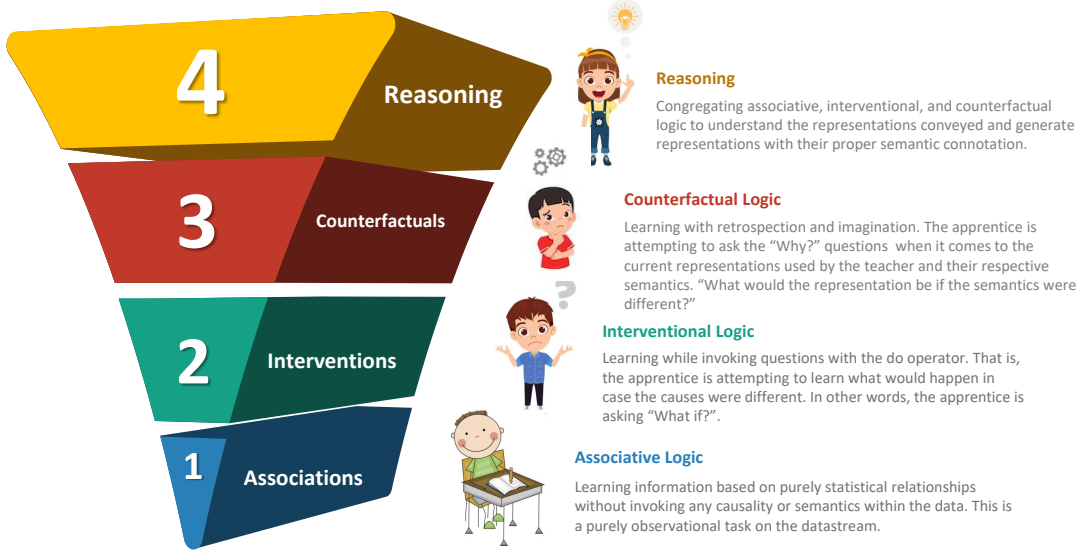


Fig. 11: Hierarchical Levels of Causality

f

specific event. In a semantic communication setting, we can re-engineer the control plane and leverage such queries (via their “do operators”) to replace simple signaling messages. As will be discussed in more detail in Section V-B, such queries can enable radio nodes to construct the causal model gradually, acquire reasoning capabilities, and ultimately learn a language.

In fact, causal inference has a long history in a variety of disciplines such as statistics, econometrics, and epidemiology [72]. Thus, various frameworks exist for studying causality. That said, existing causality frameworks cannot be applied in a plug and play fashion over the data found at the teacher and apprentice of a semantic communication system. Henceforth, it is necessary to scrutinize the key aspects of causal inference that enable building a solid causal representation learning framework for semantic communications. Next, we shed light on definitions and fundamental concepts from causal representation learning that can be leveraged to establish a strong semantic connectivity between the teacher and apprentice, and subsequently a strong semantic connectivity in an E2E network.

### B. Fundamentals of Causal Reasoning

Formalizing the concepts of causal representation learning in a semantic communication system requires introducing causality in the mechanism that builds and reasons over the semantic language and its components (semantic representations). One way to characterize such causality is via a structural causal model (SCM). Essentially, mapping our language  $\mathcal{L}$  into an SCM creates many opportunities. In fact, SCMs enable three intrinsic causal concepts: a) graphical models, b) structural equations, and c) interventional and counterfactual logic.

Henceforth, we map the semantic language between teacher and apprentice to an SCM as follows:

**Definition 5.** *Constructing a semantic language with causal reasoning capabilities requires mapping the language to an SCM  $\mathcal{L} := (\psi_L, p(\epsilon))$  where  $\psi_L = \{s_i\}_{i=1}^N$ . The learnable data can now be written:*

$$X_{l,i} := s_i(\epsilon_i, \rho_i). \quad (7)$$

Here,  $\rho_i$  is the set of direct causes leading to the specific data patterns in  $X_l$ ,  $N$  is the number of semantic content elements contained in a datastream, and  $p(\epsilon)$  is the joint distribution  $p(\epsilon) = \prod_{i=1}^N p(\epsilon_i)$  over mutually independent exogenous noisy variables. Such exogenous variables map to the variability of the data previously defined in Section IV-C. In a causal setting [74], such “noise” is thought to be an unaccounted source of variation.

Defining a language based on Definition 5 enables us to leverage queries (interventions and counterfactuals). These queries (see Fig. 7 for illustrative examples queries) enable the apprentice to understand the reason a teacher used a particular representation (to describe a specific semantic content). In principle, as seen in Fig. 11, among queries, interventions on the causality ladder rank higher than associations. We can see that, under this definition, the language is now directly related to the structural assignments rather than being concerned with mapping the datastream to its corresponding task on an associational and statistical level. Furthermore, relying on interventions enables constructing a representation that does not merely rely on the *observed datastream*. That is, the intervention proactively requires modifying the cause of events

leading to a particular datastream, which enables inferring the ultimate technique to generate semantic content. If we look back at our intuitive example in Fig. 7, we can see that acquiring understanding requires the apprentice to intervene and ask questions to ultimately build the representation. Strictly speaking, in a causal setting, an *intervention* represents a subset of the queries that the apprentice can ask to build their understanding, as defined next:

**Definition 6.** *An intervention contributing to understanding, decoding, and eventually re-generating a representation communicated by the teacher is a question posed in the form of a “do operator”. In other words, given a model mapped to a representation  $p(\mathbf{Z}|\mathbf{X}_l)$ , the apprentice ought to ask questions that enable calculating and characterizing the quantity  $p(\mathbf{Z}|do(\mathbf{X}_l = \mathbf{x}), A)$ . Here  $A$ , is a latent variable that might be affecting the outcome of the representation. The “do operator” enables the apprentice to understand the changes induced on the overall representation in case the datastream is different.*

It is important to note when intervening on a distribution with the “do operator”, we are not considering a certain sub-population for which we observe  $\mathbf{X}_l = \mathbf{x}$ , but rather we are reasoning over the changes occurring on the datastream to be conveyed after taking an action on  $\mathbf{X}_l$ , namely  $do(\mathbf{X}_l)$ . Moreover, these interventions (and the counterfactuals, which will be discussed next) can:

- 1) Replace control and signaling datastreams classically transmitted by radio nodes. This is particularly the case when the intervention is simple.
- 2) Be transmitted via classical or semantic channels, that is, they can be: a) transmitted classically by the receiver (if they lack all reasoning foundations); or b) transmitted using a semantic representation. This can asymptotically lead to reverse mentorship whereby the apprentice is teaching the teacher their semantic language and vocabulary.

Furthermore, an SCM model also enables the teacher-apprentice pair to leverage the concept of counterfactuals in a semantic communication system. Counterfactuals are at the highest level of hierarchy in causality to ultimately reach reasoning as shown in Fig. 11. That is, counterfactuals enable the apprentice to ask questions that include “why” (not only “what if”) to understand the representation, and build their knowledge base. Here, formalizing counterfactuals goes beyond the *do operator*. It incorporates factual data and an intervention (in which parts of the environment remain unchanged). To characterize the language via counterfactuals, one can modify its corresponding model as follows [74]:

**Definition 7.** *Counterfactual questions can be reflected in a language  $\mathcal{L}$  by replacing the prior distribution of variability  $p(\epsilon)$  with the posterior  $p(\epsilon|\mathbf{x})$ :*

$$\mathcal{L}_x := (\psi_L, p(\epsilon|\mathbf{x})). \quad (8)$$

In this case, mapping a language to an SCM of this form, enables the apprentice to “interrogate” their teacher with “why”s and “what if”s. The answers to these questions equip the radio node with a better understanding and knowledge base, enhancing their logic to ultimately tend to the human’s brain. This definition further showcases the benefit of building a language according to an SCMs.

So far, we have mapped a semantic language to an SCM and we have proposed a suite of queries (interventions or counterfactuals) that can operate in the control or data plane to build a common language and understanding between teacher and apprentice. Another key benefit of mapping a language to an SCM model is *the disentanglement property*. That is, the semantic representation and its respective content element can be *easily* separated from each other (the same way each word captures a standalone meaning in a natural language). Next, we highlight the disentanglement property and its respective consequences:

**Definition 8.** *Building a semantic language  $\mathcal{L}$  that can be mapped to an SCM model enables disentangling each datastream and its respective representation from other established representations. In other words, the model describing the language can be written as [74]:*

$$P(\mathbf{X}_l) = P(X_{l,1}, \dots, X_{l,N}) = \prod_{i=1}^M P(X_{l,i}|\rho_i), \quad (9)$$

where  $M \leq N$ .

Definition 8 implies that, when causality is present in the model, one can separate the parent root causes from each other, and consequently disentangle the semantic content elements. Furthermore, from Definition 8, we make the following observations:

- 1) Performing an intervention or a counterfactual on one mechanism  $P(X_{l,i}|\rho_i)$  does not change any of the other mechanisms  $P(X_{l,j}|\rho_j)$ , where  $(i \neq j)$ . This is inherently important because it creates a foundation that eases disentangling one learned task from another within a datastream at the teacher. For instance, in Fig. 7, by performing queries on Max via mathematical operators (interventions or counterfactuals), can separate Max from the background or other transmitted information.
- 2) Acquiring information about a specific mechanism  $P(X_{l,i}|\rho_i)$ , does not give us any information about  $P(X_{l,j}|\rho_j)$ . Intuitively speaking, information acquired about Max does not give further information about the sunflower in Fig. 7.

In light of this, the “independent causal mechanisms” (independent causal mechanisms might be statistically correlated) principle [73] in Definition 8 characterizes the dynamics of information shared between two distinct semantic representations (which represent two distinct semantic content elements). In essence, this is the principle that enables causal semantic representations to be invariant, autonomous, and independent

as will be discussed in Section V-C. The importance of this principle can be highlighted by the following observation: if semantic representations were instead modeled via non-causal and purely statistical techniques [72], once the apprentice poses a new query on one semantic representation, the others will also be affected. In such cases, the factorization is known to be *entangled*.

As a result, in the converse, i.e., to evaluate the causality of a language, we first define disentanglement in what follows:

**Definition 9.** For a set of representations  $\mathbf{Z}$ , s.t.  $\mathbf{X}_l = g(\mathbf{Z})$  for some mapping  $g$ , a representation is known to be *causally disentangled* if the following factorization is possible [74]:

$$p(Z_1, \dots, Z_M) = \prod_{i=1}^M p(Z_i | \rho(Z_i)), \quad (10)$$

where  $\rho(Z_i) \subset Z_{j_{i \neq j}} \cup \epsilon_i$  and  $\epsilon_i$  is the exogenous causal factor of  $Z_i$ .

From Definition 10, one can infer the following: Given a language  $\mathcal{L}_t$ , whose causality is not yet proved, if this language  $\mathcal{L}_t$  admits multiple representations that can satisfy (10), one can claim that this  $\mathcal{L}_t$  is proven to be a *causal semantic language*. Furthermore, the presence of causality in the system opens the door for utilizing causal logic and queries like counterfactuals and interventions to perform reasoning and extract further information from the exchanged and generated tasks. Meanwhile, a non-causal model would only permit inferring information in a limited i.i.d scenario.

Furthermore, one can also *leverage the principle of independent causal mechanisms to disentangle the structure and variability* of a specific datastream. In fact, one can adopt frameworks like contrastive learning [75], which is a discriminative self-ML framework that performs positive and negative sampling to ultimately create a semantic equivalence between the datastreams that need to yield the same semantic representation within a specific variability. It also establishes a distance between semantically different samples within a representation space. In fact, as shown in [76], performing contrastive learning on a causal model proved to be capable of invariantly learning representations. We will further elaborate the techniques and enablers of causal invariant representation learning next.

### C. Causality for Generalizable Representation Learning

Thus far we have scrutinized the fundamentals of causal representation learning. In particular, we have identified its peculiar features that grant the teacher and apprentice the proper tools to reach the reasoning foundation needed to extract a minimal and efficient semantic representation. That said, given that raw data can result from heterogeneous sources, and can exhibit horizontal and vertical shifts, the representation must be *invariant* to such changes. Such an invariance enables the teacher/apprentice pair to generalize the learned semantic language to *new, unseen, and out of domain, distribution, or*

*context* datastreams. Thus, this generalizability is characterized by the universality of the semantic language.

**Definition 10.** A semantic representation is dubbed, *generalizable*, if it fulfills the general causal invariant prediction criterion. That is, despite different “what if”s posed on the causal model, the same representation results in describing its respective content elements in data:

$$p^{do(\kappa_i)}(\mathbf{Y}|\mathbf{Z}) = p^{do(\kappa_j)}(\mathbf{Y}|\mathbf{Z}) \forall \kappa_i, \kappa_j \in \mathcal{K}, \quad (11)$$

where  $\mathcal{K}$  is the set of queries at the apprentice which pose interventions on the considered SCM.

Definition 10 is formulated based on the following observation: When considering a particular representation, and its corresponding content element, if a set of two different queries leads to the *exact same learned causal model*, then such a representation is *generalizable*. That is, irrespective of the queries that the apprentice has asked, the semantic language used by the apprentice remains consistent. This means that the representation used by the teacher can be applied irrespective of the data, distribution, and context. One analogy that one can draw is to “words” that we use in our daily lives: A *word* (which mimics a representation in a semantic language) consistently describes the same meaning in any context, time, and space limits. Moreover, in light of this, to guarantee that the yielded representations from our causal model are *invariant*, we further detail two approaches that leverage the invariance principle in the design of the causal reasoning faculty:

#### (i) Contrastive Causal Learning:

Given that contrastive learning is a form of self-supervised learning, one can leverage the approach adopted in [76] to identify the invariant structure properties of a semantic representation. In particular, one can train or bias the apprentice, prior to any information exchange to a neutral causal structure of the data, then the apprentice is causally taught to disentangle the structure and variability, which map to the content and style of an SCM. Thus, the structure  $\psi_L$  and the variability  $\varkappa_L$  are assumed to be independent of each other. For instance, if the structure denotes the semantic representation of a dog, that will not have a bearing on the breed of the dog. They complement each other yet remain independent. Building on this contrastive learning setting, invariance can be achieved by applying Definition 10 to the specific disentangled setting herein focused on the structure rather than the representation:

$$p^{do(\varkappa_i)}(\mathbf{Z}|\psi) = p^{do(\varkappa_j)}(\mathbf{Z}|\psi), \forall \varkappa_{i,j} \in \mathcal{V}. \quad (12)$$

That is, the knowledge base of a radio node has robustly acquired a particular structure  $\psi$  irrespective of the variability in which such structure might appear in. In other words, assuming a particular representation  $\mathbf{Z}$  describes a structure  $\psi$  mapping to a dog, and a variability mapping to  $\varkappa$  German shepherd, the radio node can robustly generalize the structure  $\psi$  to any breeds of dogs,

irrespective of where they appear (the dog example is for illustrative purposes only and this can be expanded to any data type).

(ii) **Counterfactual Invariance:**

So far, we discussed how counterfactuals can enable the apprentice to gather higher reasoning capabilities. Moreover, counterfactuals enable leveraging the framework of *counterfactual invariance* [77]. Adopting this framework can construct predictors that are invariant to particular perturbations in the raw data  $\mathbf{X}_l$ . Based on Fig. 11, such form of invariance must be stronger than the one imposed via interventions. This is important for semantic communications because for this causal model, despite different “why”s posed on the causal model (in contrast to the weaker “what if”s), the semantic language remains consistent. Essentially, this framework is built on the premise of identifying an additional variable, say  $\Upsilon$ , that captures information that must not influence the semantic representation  $\mathbf{Z}$ , nor its semantic content  $\mathbf{Y}$ . If we take the intuitive example in Fig. 7, in the process of building a semantic representation for Max, the background or the position of Max, must not affect the semantic representation chosen for Max. As such, the background does not have a causal effect on the covariates of  $\mathbf{X}_l$ . More formally, the definition of counterfactual invariance is given as follows:

**Definition 11.** *A semantic representation  $\mathbf{Z}$  is counterfactually invariant to  $Q$  if for any give sample  $\mathbf{X}_l$ , a counterfactual  $X_{l,q} \sim p(x_{l,q}|\mathbf{x}_l)$ , and  $\forall q \in \mathcal{Q}$ , we have  $z(x_{l,v}) = z(x_l)$  [77].*

Adopting this logic at the apprentice requires us to *first* identify the causal directions surrounding the raw datastream bits  $\mathbf{X}_l$  and their corresponding semantic content  $\mathbf{Y}$ . *Second*, the apprentice must be capable of capturing the attributes belonging to  $Q$ . Furthermore, the apprentice must be able to identify the associational relation between  $Q$  and  $\mathbf{Y}$  to reason whether this relationship is due to confounding or selection bias [74]. That said, acquiring all of this knowledge at this apprentice requires either a knowledge map or a set of labeled data, i.e., the raw data and their respective semantic content (not the representation). Henceforth, it is worth exploring techniques that enable leveraging this concept further while freeing the apprentice from the aforementioned restrictions.

*D. Challenges and Future Directions*

In this section we first highlight the main challenges facing causal models. Then, we discuss the potential opportunities that can be leveraged to address such challenges and subsequently build semantic communication systems with robust reasoning faculties.

- **Causality alone is not enough:** Causal models enable building a semantic language with intrinsic causality. Such causality enables the teacher/apprentice to acquire an understanding via counterfactuals and interventions. These

queries and their corresponding answers equip radio nodes with robust knowledge bases. That said, implementing SCMs into a semantic communication system faces some challenges such as a: a) difficulty in initializing such SCMs, b) difficulty in expressing causal and statistical relationships within the data. Here, an incorrect initialization might lead to a *biased* semantic language. Meanwhile, statistical relationships on top of causal ones expand the universality of a semantic language, and subsequently improve the generalizability of wireless networks. In essence, an SCM needs to be embedded into a large AI or ML whose inputs and outputs may be unstructured or naturally entangled. This also further enables expressing the statistical relationships in the data that are not characterized via causality. Here, we envision that one avenue that can potentially mitigate the aforementioned challenges is to adopt an AI system that merges connectionist AI and symbolic AI [78]. One form of this integration is dubbed, *neuro-symbolic AI*. Neuro-symbolic AI is an emerging concept that merges data-driven neural architectures which extract statistical structures from raw datastreams with symbolic AI representations of logic. In fact, in our recent work in [79], we showed that a neuro-symbolic AI framework based on generative flow networks [80] and logic-based symbolic components can be used to design an end-to-end semantic communication systems. Our results in [79] show that causality-based neuro-symbolic AI can indeed help achieve minimal representations, create symmetric communication channels, enable generalizability, and reduce the amount of data transmitted. Naturally, this early work can further be extended to accommodate several of the key concepts that we presented here, including the use of more elaborate SCM models, fully exploiting the causality ladder in Fig. 11, designing novel neuro-symbolic architectures that go beyond generative flow networks, and the use of our newly defined metrics in Section VI.

- **Highly complex causal models are problematic:** Expressing all the complexity of a task via an SCM leads to a challenge in solving the optimization problem in (4). That is, if the data is very complex, the reasoning radio node might not be able to find a tradeoff between the structure and complexity of the data. Also, this will jeopardize the explainability of the established knowledge base and the expressed semantic language. One key aspect to investigate here would be extending SCMs via novel AI techniques that can represent highly complex causal models while maintaining their expressivity.
- **How to pose the proper interventions and counterfactuals via queries?** Theoretically speaking, and given a particular datastream that exhibits structural characteristics, the apprentice needs to pose the *right* interventions and counterfactuals. In other words, the apprentice needs to have *minimal reasoning capabilities* that enable them to pose the right questions. On the one hand, if the apprentice

poses questions that are irrelevant of the context, the learned structural model will be wrong. On the other hand, if the apprentice has been accustomed to a particular context, and is therefore *biased*, they might ask questions that only enable learning selective semantic content elements of the data. Here, one aspect that is worth exploring is to send the apprentice information that enable them to build proper context and subsequently pose the proper queries. For instance, given that the teacher uses raw data to complement their semantic representations, the teacher can build on this raw data to guide the apprentice vis-à-vis context and the queries they must ask to understand the language. Here, one can adopt the game-theoretic scheme proposed in [79] to initialize the gradual design of a semantic language. Clearly, this is a nascent research direction that needs to be investigated more closely.

Thus far we have scrutinized the fundamental reasoning concepts needed to design a semantic communication system. Next, we will investigate novel metrics that enable a proper evaluation of future semantic communication systems.

## VI. SEMANTIC-ENABLED COMMUNICATION METRICS

Based on the reasoning foundations we have built so far, in this section, we will establish a set of new semantic communication metrics that enable evaluating the performance of next-generation semantic communication systems.

### A. Index of Communication Symmetry

In classical communications, the setting of communication was governed by asymmetry. This asymmetry had its bearing weight at the transmitter side, whereby transmitting nodes had the utmost power in terms of knowledge with respect to the datastream. That is, the transmitter is either generating or observing the datastream. Meanwhile, the receiver’s main goal was to only *identify*, at the destination, the message produced at the source. Thus, the receiver did not have an *active* role given that it could not *generate* anything from minimal information. This phenomenon is in alignment with the data processing inequality that explicitly states the fact that “information” (per Shannon’s definition) cannot be created in an *ex nihilo* fashion. However, in a semantic communication system, the overall situation changes given that the apprentice can leverage reasoning and causality to *generate* the originally sent message. In fact, asymptotically, one can think there could be a channel between two points, even if there is no direct physical (wired or wireless) link joining them. In other words, in semantic communication networks, the concept of communications is *governed by manipulation*. In fact, the authors in [39] claim that, “*there is no communication without manipulation*”. Thus, whatever the teacher manipulates at the transmission end, must be re-generated within its semantic context identically at the apprentice.

Consequently, it is necessary to introduce a suite of novel metrics that characterize the level of symmetry between a teacher and an apprentice. Such a metric must:

- 1) Understand the reasoning capability of the teacher and the apprentice.
- 2) Investigate whether an equilibrium has been reached, whereby the teacher and apprentice can majorly rely on semantic-based transmissions.
- 3) Scrutinize whether a reverse-mentorship situation is taking place, whereby the apprentice’s reasoning capabilities are superior to the teacher’s.

Prior to proposing such a metric, we first need to define a novel concept, called *semantic impact*. This stems from the fact that the significance of a particular semantic representation  $Z_i$  can be characterized by the equivalent number of data packets one would have needed to convey the exact same message. Considering a particular semantic representation,  $Z_i$ , and its respective semantic content element  $Y_i$ , one must ask: “If we transmitted the data classically, during a time duration  $\tau$ , how many data packets would be needed to generate the same semantic content?”. We answer this question via a novel metric, dubbed, semantic impact:

**Definition 12.** *The semantic impact  $\iota_\tau$  generated by a semantic representation  $Z_i$  during a time  $\tau$  is the number of packets that would have been needed to be transmitted to regenerate the semantic content element  $Y_i$ .*

**Proposition 2.** *The communication symmetry index between a teacher  $b$  and apprentice  $d$ , for a transmission session  $\tau$  is given by:*

$$\eta_{b,d,\tau} = \frac{\zeta_{d,\tau}}{\nu_{b,\tau}} \times \iota_{\tau,Y_i}, \quad (13)$$

where  $\iota_\tau$  is the semantic impact,  $\zeta_{d,\tau}$  is the number of query packets (interventions, counterfactuals, etc.) requested by the apprentice to reason over the message transmitted, and  $\nu_{b,\tau}$  is the number of raw data packets transmitted by the teacher to accompany the semantic representation sent. We also note that, if the queries are communicated via a semantic representation to the teacher, one can find  $\zeta_{d,\tau}$  by applying the concept of semantic impact on the representation used.

Thus, this communication symmetry index along with the semantic impact enable us to characterize the reasoning state of the teacher and apprentice and the equilibrium reached by the teacher and the apprentice:

- If  $\eta_{b,d,\tau} \leq 1$  and  $\iota_\tau > 1$ : This is a setting in which the apprentice has little to no knowledge base. Here, the apprentice has not acquired the reasoning capacity to intervene and interrogate the teacher on the semantic representation sent. This setting asymptotically mimics the pure classical communication scenario, whereby the majority of the data is still being sent in its raw form to complement the semantic representation.
- If  $\eta_{b,d,\tau} \rightarrow \iota_\tau$  and  $\iota_\tau > 1$ : This is a setting in which the apprentice has considerable knowledge and reasoning faculties. The teacher is still complementing their conveyed information with raw data, however the apprentice

is intervening often to understand the causal structure of the data and gradually rely on semantic representations.

- If  $\eta_{b,d,\tau} = \iota_\tau$  and  $\iota_\tau > 1$ : In this setting, the apprentice is intervening the same way current receivers transmit an acknowledgment, meanwhile the teacher is only relying on raw data transmissions to characterize the unlearnable part of the data (that is solely governed by randomness).
- If  $\eta_{b,d,\tau} > \iota_\tau$  and  $\iota_\tau > 1$ : This is a very peculiar setting, whereby the datastream is mostly learnable and not memorizable, thus the teacher is relying mostly on semantic representations. Also, the apprentice is actively intervening to generate the conveyed message from the set of semantic representations conveyed.
- If  $\eta_{b,d,\tau} > \iota_\tau$  and  $\iota_\tau \leq 1$ : Here, the number of queries requested by the apprentice is larger than raw data transmissions by the teacher. This settings represents a challenging scenario in which the teacher is unable to extract a proper semantic representation to be communicated with the apprentice.

We can see that parameter  $\eta_{b,d,\tau}$  does not go significantly below 1 unless a defect in reasoning is observed. Next, we build on the concept of communication symmetry index to propose a novel *reasoning capacity* metric.

### B. From Information Capacity to Reasoning Capacity

In classical communications, most of the network metrics are built on the concept of information capacity. In essence, information capacity characterizes the maximum achievable data rate of a particular communication link, under specific bandwidth allocation. This metric is fundamentally important because it enables the system designer to quantify the range of communication rates possible, and thus understand the type of applications or use cases that can benefit from such a QoS. That said, it is important to note that this characterization is based on the classical concept of “information”, whereby information is merely quantifying uncertainty. In other words, the capacity will not describe the number of symbols that can be contained in a communication channel, but rather the number of choices that the “*information*” can take, i.e., the number of choices one could transmit in a particular communication channel. Shannon proved that the channel capacity is equal to the mutual information of the channel maximized over all possible choices at the transmitter:

$$C = \max_{P(X)} I(X, \tilde{X}), \quad (14)$$

where  $\tilde{X}$  is the output of the classical communication channel. This definition is built on the method Shannon utilized to quantify information via uncertainty (see Section III-B for a more elaborate discussion about information via uncertainty). Building on this definition, there exists a need to propose a novel “capacity” metric that is capable of accounting to the novel nature of a semantic communication system. To do so, the following key aspects must be scrutinized prior to the design of a novel notion of *reasoning capacity*:

- 1) Information must be defined as a semantic substance rather than a mere uncertainty measure [81], whereby the quantification of “information”, characterizes the semantic representation communicated between the teacher and the apprentice.
- 2) The concept of “information transmission” must not be viewed as a mere propagation between two points of physical space [39]. Instead, it must be seen as a generative process at the apprentice that requires reasoning, and is controlled via manipulation.
- 3) Reasoning capacity cannot be high in the presence of an asymmetry between the teacher and the apprentice. In a classical setting, the communication is highly manipulated by the transmitter and the receiver does not perform any tasks on the data beside mere decoding and reconstructing. Furthermore, in a semantic setting, having a very well-versed teacher but a non-trained apprentice cannot yield a high reasoning capacity, as the approach ends up leading to a scenario similar to the classical asymmetrical one. That said, in contrast to the classical approach, the apprentice is in the gradual reasoning stage as they are incrementally building on their knowledge.
- 4) Similarly, having a well-versed apprentice but a teacher with weak representational capability leads to the use of redundant computing and communication resources over the course of connectivity. Even though the main course of information transmission must flow from the teacher to the apprentice, the teacher must gradually acquire new tools from the apprentice. That is, the teacher will establish a subset of their semantic language based on the apprentice’s previous experiences and knowledge base.
- 5) Henceforth, a high reasoning capacity necessitates a teacher with well-versed representational skills (one that can easily build a semantic language based on a particular datastream), and an apprentice with good reasoning capabilities (one that can understand a representation and use it to computationally generate semantic content).

Our novel, proposed reasoning capacity measure must, unlike the classical capacity measure, correlate with the message complexity proposed in Proposition 4. That is, the classical capacity measures would not change as a function of the *message complexity* but rather as a function of the *channel conditions*. If the goal is to load a large XR content, a high data would be needed as the information is recovered at the receiver side in a *bit-by-bit* fashion. Instead, if the receiver becomes an apprentice that relies on learning the content rather than simply recovering it, the *understanding* of the apprentice and its *impact* on the reconstruction process would be the KPI of a reliable communication link between teacher and apprentice. Essentially, characterizing the understanding of the apprentice depends on:

- 1) Reasoning as measured by the number of queries posed by the apprentice.
- 2) Efficiency and minimalism as measured by the number of raw messages supplemented to the semantic representa-

tion, as well as the impact of the semantic representation. Consequently, we propose the following KPI for semantic communication links:

**Proposition 3.** *The reasoning capacity between a teacher  $b$  and an apprentice  $d$  is given by:*

$$C_R = \Omega \log_2(1 + \eta_{b,d}), \quad (15)$$

where  $\Omega$  is the maximum computing capability of the server used to represent or generate the semantic representation, and  $\eta_{b,d}$  is the communication symmetry index per second.

Here,  $\eta_{b,d}$  is computed per second rather than transmission session, thus  $\tau$  is integrated out. Furthermore, the reasoning capacity is still a binary log function since the representation is still sent via bits over a channel. This proposition is universal in a sense that it is independent of the type of semantic representation utilized. It is also additive to the classical capacity metric given that  $\iota_{\tau, Y_i}$  enables the classical conversion. It is thus important to find techniques that can accurately measure  $\iota_{\tau, Y_i}$ . Furthermore, given a datastream  $X$ , with  $\mathbf{X}_l$  learnable information and  $\mathbf{X}_m$  memorizable ones, one can simply use addition to compute the capacity. That is, the total achievable capacity would be given by:

$$C_T = C_C + C_R = W \log_2(1 + \gamma) + \Omega \log_2(1 + \eta_{b,d}), \quad (16)$$

where  $W$  is the bandwidth and  $\gamma$  is the signal-to-interference-plus-noise-ratio (SINR). One important thing to note here is that while the  $C_C$  is limited by Shannon's bound, and the bandwidth of operation;  $C_R$  is essentially bounded by the reasoning bound and the computing resources. Notably, claiming an E2E capacity that might achieve an impact that is higher than what Shannon could have is possible because we are alternatively relying on computing. In other words, if the apprentice is capable of regenerating the message (via a semantic representation) faster than it could be classically transmitted, the E2E capacity could asymptotically reach beyond Shannon's limit.

Hence, we have so far elaborately investigated the principles of reasoning to elicit a semantic representation. Subsequently, we have proposed a set of novel semantic-based communication metrics to ensure that future systems are evaluated efficiently. Next, we will look into the considerations needed when scaling the semantic communication system from one teacher and apprentice, to an overall large scale network.

## VII. SCALING SEMANTIC COMMUNICATIONS: FROM SEMANTIC LINKS TO SEMANTIC NETWORKS (6G AND BEYOND)

Thus far, we have discussed the concept of semantic communication systems for a link a teacher and an apprentice, i.e., two radio nodes. In this section, we further discuss how those can be used to build large scale semantic communication networks.

### A. Early deployments of semantic communications: challenges and facilitators

In this section, we will first discuss the techniques that enable developing the reasoning skills and the language with nascent and uninformed teachers. Then, we will shed light on the early use-cases and forms of semantic communication systems.

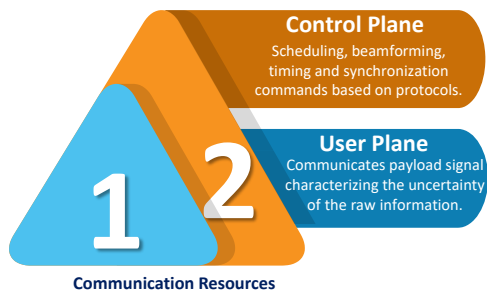
1) *Uninformed Teachers:* A semantic communication system essentially depends on the reasoning faculties of the teacher and apprentice. While the symmetry between the teacher/apprentice pair increases as the apprentice learns the semantic language, we have assumed thus far that the teacher has acquired the faculty to provide and teach a semantic language. However, in many instances, the teacher lacks the capability to develop their reasoning to ultimately teach the apprentice this language. This is the case when a radio node has never communicated via semantic-language before, or when this node is trying to communicate a novel structure that is considered "novel vocabulary" with respect to its acquired knowledge. In this case, this "uninformed" apprentice must use alternative techniques to develop their reasoning capability and learn the foundations of a language, given a structural datastream. This can be performed via:

- **Reverse mentorship:** If the apprentice has powerful reasoning and representational capabilities, the apprentice can through, interventions and counterfactuals, teach the teacher specific data representations. The teacher may in fact "meet" such an informed apprentice if it moves in a network.
- **Data showers:** The teacher can resort to standalone cloud services that offload standalone libraries that can complement their knowledge and ultimately build a basic semantic language.

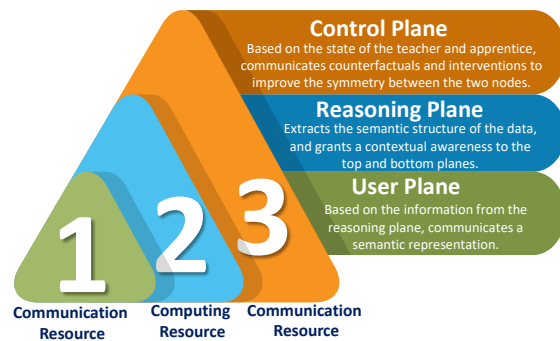
2) *How to prepare today's networks so as to deploy semantic communications?:* Thus far, we have investigated how the concept of semantic communications can pave transform wireless networks from AI-augmented, data driven networks, to AI-native reasoning-driven networks. While this leap can be revolutionary, moving towards this direction will require:

- **Readiness of computing resources:** In semantic communication systems, producing/understanding a language on the transmitter/receiver end requires a high abundance of computing resources at the end devices and UEs. In light of the AI-wave for wireless networks, an exponential growth has been observed in the computing resources of modern devices, yet, this remains heterogeneous and disparate among devices. Henceforth, an open problem in this direction is to scrutinize: a) the level of reasoning that can be achieved with limited computing resources, and b) the evolution of a semantic language built based on this reasoning, and whether it can yield substantial benefits to the E2E system. In addition, it would be important to investigate how a semantic communication network can be built out of devices that have very





5G and Beyond Communication Planes



6G and Beyond AI-Native Planes

Fig. 12: An illustrative figure showcasing the transition from communication-only planes in 5G and beyond to AI-Native Planes in 6G and beyond.

different reasoning/computing capabilities, and what the performance gains would be for semantic communication in such heterogeneous systems.

- Maturity of reasoning-intensive applications:** As a result of the characteristics of semantic communication systems, some applications might benefit more than others from using semantic communication mechanisms. For instance, services that require a high level of service intelligence (in contrast to operational intelligence which is used to optimize the network to fulfill the QoS of the application), i.e., the execution of the service is complex and/or might have a high level of autonomy; can highly benefit from the semantic structure to produce content. For example, the production of holograms in holographic teleportation and high-precision manufacturing in Industry 5.0 are tasks that require multiple AI mechanisms (this is independent of how the content will/can be transmitted). Here, using a semantic language to communicate the semantic structure modifies the content production process at the receiver. That is, instead of decoding bit-by-bit such complex messages, the receiver must reproduce (using their computing resources) the described task. Instead of transmitting the complex content in a message, the message would be comprised of a language that describes the key aspects of the complex content. This enables teaching the receiver to *automate and produce the content based on their knowledge base*.
- Highly cooperative systems:** Semantic communication systems highly depends on the concept of a *language*. Services that require continuous and real-time data or control messaging between homogeneous nodes can benefit from a language more than others. That is, a swarm of vehicles or robots interacting with an environment can highly benefit from communicating a language. This can minimize the

number of control messages repetitively communicated to achieve a goal. Also, digital twins [42] require a tight level of synchronization between the physical and the cyber twin. A semantic language here can improve the cooperation between the twins and guarantee achieving a real-time replica of twin in the metaverse.

- Evolution of AI:** On top of the evolution needed on the network's architecture, semantic communication requires a major evolution of current AI frameworks. As of today, causal representation learning, as well as frameworks that deploy associational and causal logic remain in their infancy. Thus, a suite of novel AI mechanisms that can efficiently execute the causal ladder we previously described is needed to ensure the deployment of future semantic networks.

Hence, based on our previous observations, we expect that the use case of *goal-oriented communications* will be the first major application of semantic communications, which somewhat explains the high interest that it received from early-on semantic communication papers [28], [29]. Here, services that could highly benefit from goal-oriented communications like cyber-physical systems, digital twins, and connected robotics and autonomous systems (CRAS) could be the first adopters of semantic communications given the: a) The abundance of computing resources as a result of their infrastructure, and b) The high level of cooperation that intrinsically exists in the operation of such services. For instance, for digital twins, transforming the communication between the physical twin and the cyber twin to a semantic-based one enables them reaching a high level of synchronization. Here, this can be done by majorly relying on computing resources, and thus *minimizing the need* for a highly reliable and low latency channel to achieve the overarching synchronization goal.

Next, we investigate the opportunities and challenges of

semantic communications in the large-scale cellular network context.

### B. E2E Semantic Large-Scale Wireless Networks

In this subsection, we further expand our study on semantic communication networks, by examining how their corresponding networking protocols could look like. Then, we investigate the techniques needed to distribute computing resources in a network, which showcases the true convergence of communication and communication resources in semantic systems. Finally, we scrutinize how semantic communication systems transform today's 5G O-RAN.

#### 1) Networking and Computing Considerations:

• **On Semantic-Based Networking:** In Section V, we investigated how interventions and counterfactuals enable the apprentice to learn about the structural causal model that the teacher is trying to convey. Mathematically, such queries are PDFs with “do operators” that enable the apprentice to better understand the causal model. That said, the apprentice will be communicating such “do operators” in the uplink in contrast to classical control acknowledgements and non-acknowledgements. In essence, classically, the “networking task” was to signal the success or failure in the reception of a particular message (or a part of a message). Meanwhile, in semantic communication systems, the mere reception of a message does not have the same significance anymore. This is why, interventions and counterfactuals should be key features of the semantic-based control messages. That is, the goal of semantic networking relies on: a) confirming the “understanding” of a semantic representation, and b) the ability to build or compile a semantic language with a similar view-point as the teacher.

Moreover, in current 5G cellular systems, the control and user planes are separated as shown in Fig. 12. This has been designed to enable a high flexibility and interoperability. On the one hand, the control plane can interact with multiple user planes. On the other hand, the user plane function can be shared by multiple control plane functions. This separation is also aligned with the O-RAN's initiative to have an *open, intelligent, virtualized and fully interoperable RAN* [82]. In contrast, in semantic cellular systems, and given the full seamless convergence of computing and communication resources, a novel *reasoning* plane is in operation. This reasoning plane is separated from the other planes, yet is *sandwiched* between the control and user plane. The main functions of this reasoning plane are centered around:

- 1) Deliver the user plane the extracted semantic representations that correspond to the source message.
- 2) Tune the control messages so as to enable building a common semantic language between the teacher and the apprentice, and asymptotically reaching a symmetry between the teacher and the apprentice.

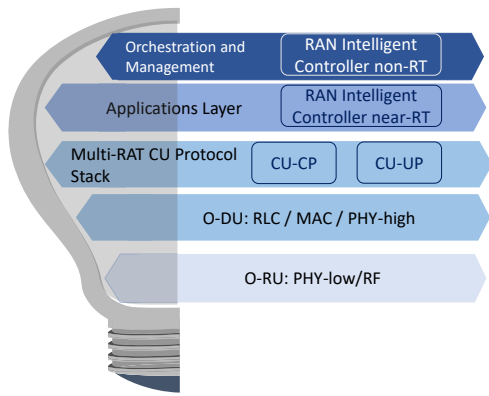
In other words, the reasoning plane will act as a master to guide the user and control plane. The operation of these two planes would be further reinforced with intent and enable reaching the overall system goals more efficiently. Based on these insights

we can make the following observations on the evolution of the physical layer and networking layer in semantic networks:

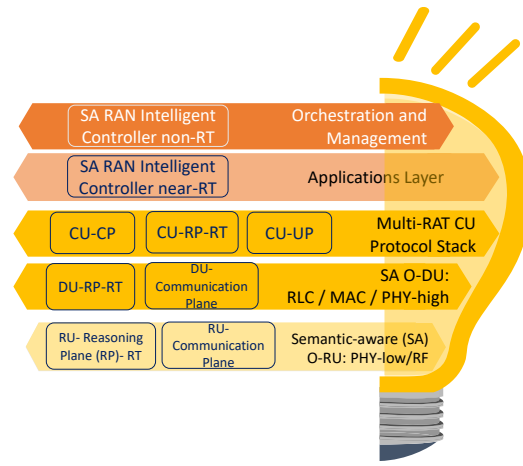
- *On the Physical Layer:* We can clearly see that semantic communications is not a substitute physical layer. In essence, semantic communications re-engineers the physical layer to be viewed as the language medium, rather than the bit-to-symbol mappers.
- *On Data-link and Networking Functions:* Networking functions like multiple access, multiplexing, resource allocation, and scheduling can conserve its current classical structure in semantic networks. Nonetheless, to improve the efficiency of the overall system, one can implement them via queries (interventions or counterfactuals) or transmit them via a semantic language. That is, if end nodes have reached a higher level of maturity in terms of the language used, the reasoning plane will require the user plane and control plane the semantic representations *that describe the next multiple access command*. This also improves the radio nodes intelligence and autonomy in making control and networking decisions. Moreover, novel semantic based networking functions can enhance the performance of time-critical communication. For instance, association, beamforming, and scheduling can be tuned based on the needs of the users as well as the changes in the environment that can be concluded from the radio nodes' knowledge bases. Essentially, all of these protocols can now be revisited with the lense of the presence of semantics, reasoning faculties, and other intrinsic features of a reasoning-driven, AI-native semantic wireless system.

• **On Reasoning and Computing-based Optimization:** In Section VII-A, we have highlighted that the paucity of computing resources may restrain the potential for future semantic communication networks. In fact, similarly to the chase of bandwidth that was observed in the migration from 4G to 5G, the evolution of semantic communications will be governed by a *chase behind computing*. Hence, optimizing computing resources, orchestrating them, and distributing them in an efficient way is an important practice of future semantic communications in 6G and beyond. The following key open-problems can be investigated:

- *Ubiquitous Distributed Computing:* In case the network has a high heterogeneity of radio nodes (with respect to the computing capability, e.g., IoT sensor communicating with digital twin), one can consider distributing the computing resources. That is, to ensure that reasoning-based services can benefit from a semantic language, one can distribute the computing resources around local low-power devices. For example, if a device cannot be augmented with more computing, one can move part of the reasoning mechanism to the edge (of the associated access point or base station, instead of the end-device). Hence, from the IoT sensor's perspective, the communication is being



5G and Beyond O-RAN Architecture



6G and Beyond AI-Native O-RAN Architecture

Fig. 13: An illustrative figure showcasing the transition of the current O-RAN architecture to a future AI-native O-RAN architecture.

carried out classically, since that link has remained intact. Meanwhile, from the edge, all the way to the digital twin, a semantic-based communication scheme is performed on that end of the link. This enables achieving fairness in the network despite the existing heterogeneity in UEs. Clearly, optimizing the distribution of computing resources for semantic communication systems is an open problem that can be carried out using distributed AI mechanisms or game theory [83].

- *Competitive and Cooperative Languages*: Throughout this tutorial, we have carried out the task of “building and learning” the language from the teacher and apprentice’s side in a fully *cooperative* way. Nonetheless, when the setting expands to multiple radio nodes, the goal of each might be distinct. Consequently, this might result in a noncooperative mechanism in building the language. Also, a subset of the radio nodes might share some common system goals, and thus might have consensus on part of the language. To examine all of these considerations, one can resort to game theory [84], in general, and *hypergame theory* [85], in particular. Essentially, hypergame theory is the confluence of game theory and decision theory, and it provides a set of tools that can be used to characterize the interaction between different radio nodes so as to model: a) The final “common” language built with respect to all nodes, b) The exact and specific view-point of every radio node with respect to the language, and c) The evolution of the “emergent” language versus the system goals of each node. Here, many other game-theoretic tools can also be explored, such as the use of signaling games and their extension (see our key results in this space [79]), referential games, and even a more complex mix of cooperative and noncooperative games. Indeed, this is an essential area of research to guarantee the

successful evolution of an emergent language, as well as the scalability of semantic communications over a large scale network.

2) *Semantic Communications in ORAN*: Recently, there has been a concerted global effort from academia, industry, and governmental agencies towards the principles of *openness* and *intelligence* [86]. As a result, the O-RAN Alliance was formally defined to reach these goals. The main objective of this alliance is to move cellular networks architectures towards disaggregated, intelligent, virtualized, and fully interoperable RAN. We can see on the left hand side of Fig. 13 the current proposed architecture of O-RAN for 5G systems. Remarkably, the O-RAN architecture is notorious for its RAN intelligent controller (RIC) which consists of two main units [87]:

- *The non-real-time (RT)-RIC* which operates in the orchestration and management plane and is used to run AI tasks that can tolerate a considerable execution time. For example, network benchmarking, AI/ML lifecycle management, and orchestration.
- *The near-RT-RIC* which operates in the applications layer and can be used for tasks that require a faster decision making (e.g. handover decisions, QoS control, and load balancing).

In fact, the “virtualized and intelligent” aspects of O-RAN have already been poised for a convergence in the computing and communication resources, however this was observed in the upper layers: applications layer and orchestration and management layer. With semantic communications, this convergence will become seamless and take place in every layer of the network stack. As a result of the introduction of the reasoning plane, we can see in Fig. 13 that novel *real-time* AI-oriented blocks have been proposed in the open radio unit (O-RU), open distributed unit (O-DU), and centralized unit (CU) layers. Moreover, we can see three planes on the CU layer, this

is where the majority of the interactions take place between the control, reasoning, and user plane we have previously highlighted in Fig. 12. We can further make the following observations with regards to the 6G and Beyond AI-Native O-RAN architecture:

- Semantic communications not only is the path for AI-nativeness, but it is also a driving force that aligns with the initiative of O-RAN. That is, with semantic communications, the intelligence chased will become operable in every aspect of the network and in three time-scales: a) real-time, b) near-real time, and c) non-real time.
- The reasoning plane introduced by semantic communications will slightly modify the non-RT RIC and RT-RIC functions, whereby their operation will be more intent-based. This intent will be formulated as a feedback from the reasoning plane and results from: a) the context of information, b) the overall system goal. For instance, the xApps (applications running on the near-RT RIC) and rApps (applications running on the non-RT RIC) can benefit from the semantic content elements identified to make better AI decisions.
- Manual requests that are still driven by the network operator can be automated and inferred from the identified root-causes of semantic content elements.

Henceforth, in this section, we have investigated the key techniques that enable a nascent teachers to get accustomed to a language, then we have highlighted the use-cases that will prospectively be the early adopters of semantic communications, as well as the key facilitators that enable expediting the maturity of semantic systems. Subsequently, we have investigated the networking and computing considerations that must take place to expand semantic communications to large scale networks. Finally, we have emphasized the role of semantic communications in O-RAN.

## VIII. CONCLUSIONS AND RECOMMENDATIONS

Semantic communications can potentially revolutionize the wireless industry, and provide a fundamentally novel way to design and operate communication systems. However, as discussed in this tutorial, in order to reap the real benefits of semantic communications, it is necessary not to reduce this area into yet another incremental extension of existing techniques such as source coding, data compression, application-aware scheduling, and natural language processing. In contrast, we advocate for creating new, rigorous mathematical foundations for semantic communication systems that lie at the intersection of AI, communication theory, networking, causal reasoning, information theory, transfer learning, and minimum description length theory. In particular, we have identified the five main tenet of semantic communication systems that include: a) minimally sufficient representations, b) semantic language, c) reasoning via causality, d) semantic-based KPIs, and e) judicious use of computing resources. We have then developed a comprehensive roadmap towards designing these pillars and building next-generation semantic communication

networks that are grounded on rigorous, concrete, and flexible knowledge-driven AI frameworks. In doing so, we have also showed that there is a need to revisit the fundamentals of information theory, and extend it to a semantic information theory that hinges on a *minimal, generalizable, and efficient* semantic language which ensures a symmetrical communication. Through the proposed frameworks, concepts, and vision, this tutorial provides, for the first time, holistic and technically-grounded answers to the following key questions:

- What is a semantic communication system, and how is it different from what we already know?
- What are the fundamental building blocks of a semantic communication system?
- How do we build the reasoning faculty of a semantic system, and how does it communicate via a *minimal, generalizable, and efficient* semantic language?
- How do we evaluate the performance of semantic communication networks, and what are the major influencers of the performance compared to classical communications?
- How do we expand semantic communications to existing and future large scale wireless networks?

In light of our panoramic investigation, we conclude with several *recommendations* to ensure a proper deployment of future semantic communication networks:

- 1) **Semantic communications is beyond goal-oriented communications:** We acknowledge that goal-oriented communication is fundamental, particularly for use-cases with common radio nodes attempting to achieve a common goal. Also, we acknowledge that such use-cases might be the first to benefit from the concept of semantic communications. Nonetheless, semantic communications is beyond goal-oriented communications and is the path to creating a fundamentally novel type of networks that we called reasoning-driven AI-native wireless networks. These new reasoning-driven, AI-native systems will be able to cater to the complex requirements of services like the metaverse, XR, and the internet of senses.
- 2) **Advances in AI and computing:** Indeed there are many wireless communication challenges in deploying semantic communications. Nonetheless, given that reasoning is the central pillar of semantic communications, various advances and developments must occur in AI so as to develop radio nodes that can build comprehensive and organized knowledge bases. Also, further computing advances are needed so that logical conclusions performed by radio nodes can meet the time-critical needs of beyond 6G applications.
- 3) **On the relationship between semantic communication systems and classical communications:** In many recent works, the concept of semantic communication systems has been touted as the ultimate and only replacement of classical communication systems, and the solution to every wireless communication challenge of the next decade. While we agree that reasoning-driven, AI-native

semantic communication systems, if built correctly, can fundamentally change the way we design wireless networks, the reality is that semantic networks and classical networks will have to co-exist and work hand-in-hand. As we outlined in Section IV, raw datastream is not entirely structured into semantic content elements and contains a lot of random information. Such random information are *better memorized than learned*, and thus must be transmitted via classical communication channels. Meanwhile, learnable data will be the key input that will be transformed into a semantic language. Therefore, we recommend that research in this space heed and acknowledge the differences between a) the important short-term needs of wireless systems (e.g., better managing mmWave and THz links, enhancing classical reliability, taming the E2E communication latency), b) a medium-term milestone during which both classical and reasoning-driven semantic networks will harmoniously coexist, serving different applications and use cases, and c) the longer-term vision of pure reasoning-driven AI systems in which the majority of radio nodes will be able to leverage their accumulated and organized knowledge base to perform versatile and logical decisions across the networking stack. For the next decade, research along all three lines must concurrently take place in order for us to truly usher in a revolution in wireless systems

- 4) **Less spectrum use through convergence of computing and communication:** The deployment of semantic communications, which will be crowning of the convergence of communications and computing, will help alleviate the technical and regulatory burdens associated with the need to open up new spectrum bands each time a new wireless cellular system generation must be deployed or a new use case of the spectrum emerges. Therefore, it is necessary for the community to further think about this convergence and its implications to current and future spectrum-related roadmaps and challenges.

In a nutshell, by answering these questions and concretely laying the foundations of semantic communication networks through a unified and systematic treatment of the underlying challenges, this tutorial is poised to become a primary reference in this burgeoning field.

#### REFERENCES

- [1] C. Chaccour, M. N. Soorki, W. Saad, M. Bennis, P. Popovski, and M. Debbah, "Seven defining features of terahertz (THz) wireless systems: A fellowship of communication and sensing," *IEEE Communications Surveys & Tutorials*, vol. 24, no. 2, pp. 967–993, Jan. 2022.
- [2] M. Soltani, V. Pourahmadi, A. Mirzaei, and H. Sheikhzadeh, "Deep learning-based channel estimation," *IEEE Communications Letters*, vol. 23, no. 4, pp. 652–655, Feb. 2019.
- [3] N. Shlezinger, N. Farsad, Y. C. Eldar, and A. J. Goldsmith, "Viterbinet: A deep learning based viterbi algorithm for symbol detection," *IEEE Transactions on Wireless Communications*, vol. 19, no. 5, pp. 3319–3331, Feb. 2020.
- [4] T. S. Cousik, V. K. Shah, J. H. Reed, T. Erpek, and Y. E. Sagduyu, "Fast initial access with deep learning for beam prediction in 5g mmwave networks," in *Proc. of IEEE Military Communications Conference (MILCOM)*, San Diego, CA, USA, Dec. 2021, pp. 664–669.
- [5] X. Li, X. Hu, R. Zhang, and L. Yang, "Routing protocol design for underwater optical wireless sensor networks: A multiagent reinforcement learning approach," *IEEE Internet of Things Journal*, vol. 7, no. 10, pp. 9805–9818, Apr. 2020.
- [6] S. Jayaprakash, M. D. Nagarajan, R. P. d. Prado, S. Subramanian, and P. B. Divakarachari, "A systematic review of energy management strategies for resource allocation in the cloud: Clustering, optimization and machine learning," *Energies*, vol. 14, no. 17, p. 5322, Aug. 2021.
- [7] S. Ayoubi, N. Limam, M. A. Salahuddin, N. Shahriar, R. Boutaba, F. Estrada-Solano, and O. M. Caicedo, "Machine learning for cognitive network management," *IEEE Communications Magazine*, vol. 56, no. 1, pp. 158–165, Jan. 2018.
- [8] Y. C. Eldar, A. Goldsmith, D. Gündüz, and H. V. Poor, *Machine Learning and Wireless Communications*. Cambridge University Press, 2022.
- [9] J. Hoydis, F. A. Aoudia, A. Valcarce, and H. Viswanathan, "Toward a 6G AI-native air interface," *IEEE Communications Magazine*, vol. 59, no. 5, pp. 76–81, May 2021.
- [10] W. Wu, C. Zhou, M. Li, H. Wu, H. Zhou, N. Zhang, X. S. Shen, and W. Zhuang, "AI-native network slicing for 6g networks," *IEEE Wireless Communications*, vol. 29, no. 1, pp. 96–103, Feb 2022.
- [11] P. Carbone, G. Dan, J. Gross, B. Goeransson, and M. Petrova, "Neuroran: rethinking virtualization for ai-native radio access networks in 6g," *arXiv preprint arXiv:2104.08111*, 2021.
- [12] T. O'Shea and J. Hoydis, "An introduction to deep learning for the physical layer," *IEEE Transactions on Cognitive Communications and Networking*, vol. 3, no. 4, pp. 563–575, 2017.
- [13] Q. Lan, D. Wen, Z. Zhang, Q. Zeng, X. Chen, P. Popovski, and K. Huang, "What is semantic communication? a view on conveying meaning in the era of machine intelligence," *Journal of Communications and Information Networks*, vol. 6, no. 4, pp. 336–371, Dec. 2021.
- [14] W. Weaver, "Recent contributions to the mathematical theory of communication," *ETC: a review of general semantics*, pp. 261–281, 1953.
- [15] Q. Lan, D. Wen, Z. Zhang, Q. Zeng, X. Chen, P. Popovski, and K. Huang, "What is semantic communication? a view on conveying meaning in the era of machine intelligence," *Journal of Communications and Information Networks*, vol. 6, no. 4, pp. 336–371, Dec. 2021.
- [16] Z. Qin, X. Tao, J. Lu, and G. Y. Li, "Semantic communications: Principles and challenges," *arXiv preprint arXiv:2201.01389*, 2021.
- [17] M. Kalfa, M. Gok, A. Atalik, B. Tegin, T. M. Duman, and O. Arıkan, "Towards goal-oriented semantic signal processing: Applications and future challenges," *Digital Signal Processing*, vol. 119, p. 103134, Dec. 2021.
- [18] E. C. Strinati and S. Barbarossa, "6G networks: Beyond shannon towards semantic and goal-oriented communications," *Computer Networks*, vol. 190, p. 107930, 2021.
- [19] D. Gunduz, Z. Qin, I. E. Aguerri, H. S. Dhillon, Z. Yang, A. Yener, K. K. Wong, and C.-B. Chae, "Beyond transmitting bits: Context, semantics, and task-oriented communications," *arXiv preprint arXiv:2207.09353*, 2022.
- [20] X. Luo, H.-H. Chen, and Q. Guo, "Semantic communications: Overview, open issues, and future research directions," *IEEE Wireless Communications*, vol. 29, no. 1, pp. 210–219, Jan. 2022.
- [21] K. Niu, J. Dai, S. Yao, S. Wang, Z. Si, X. Qin, and P. Zhang, "Towards semantic communications: A paradigm shift," *arXiv preprint arXiv:2203.06692*, 2022.
- [22] W. Yang, H. Du, Z. Liew, W. Y. B. Lim, Z. Xiong, D. Niyato, X. Chi, X. S. Shen, and C. Miao, "Semantic communications for 6G future internet: Fundamentals, applications, and challenges," *arXiv preprint arXiv:2207.00427*, 2022.
- [23] C. E. Shannon, "A mathematical theory of communication," *The Bell system technical journal*, vol. 27, no. 3, pp. 379–423, Jul. 1948.
- [24] Y. Zhang and D. A. Adjeroh, "Prediction by partial approximate matching for lossless image compression," *IEEE transactions on image processing*, vol. 17, no. 6, pp. 924–935, Apr. 2008.
- [25] D. Sculley and C. E. Brodley, "Compression and machine learning: A new perspective on feature space vectors," in *Proc. of Data Compression Conference (DCC'06)*, Snowbird, UT, USA, Mar. 2006, pp. 332–341.
- [26] H. Xie, Z. Qin, X. Tao, and K. B. Letaief, "Task-oriented multi-user semantic communications," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 9, pp. 2584–2597, Jul. 2022.
- [27] F. Binucci, P. Banelli, P. Di Lorenzo, and S. Barbarossa, "Dynamic resource allocation for multi-user goal-oriented communications at the wireless edge," in *Proc. of the 30th European Signal Processing Conference (EUSIPCO)*. Belgrade, Serbia: IEEE, Aug. 2022, pp. 697–701.
- [28] M. K. Farshbafan, W. Saad, and M. Debbah, "Curriculum learning for goal-oriented semantic communications with a common language," *arXiv preprint arXiv:2204.10429*, 2022.
- [29] C. Zhang, H. Zou, S. Lasaulce, W. Saad, M. Kountouris, and M. Bennis, "Goal-oriented communications for the iot and application to data compression," *IEEE Internet of Things Magazine, Special Issue on Pervasive, Efficient and Smart Signal Processing for IoT*, to appear 2022.
- [30] H. Shajaiah, M. Ghorbanzadeh, A. Abdelhadi, and C. Clancy, "Application-aware resource allocation based on channel information for cellular networks," in *Proc. of IEEE Wireless Communications and Networking Conference (WCNC)*, Marrakesh, Morocco, Apr. 2019, pp. 1–6.
- [31] M. Chen, W. Saad, C. Yin, and M. Debbah, "Data correlation-aware resource management in wireless virtual reality (VR): An echo state transfer learning approach," *IEEE Transactions on Communications*, vol. 67, no. 6, pp. 4267–4280, Feb. 2019.
- [32] P. Temdee and R. Prasad, *Context-aware communication and computing: Applications for smart environment*. Springer, 2018.
- [33] P. E. G. Silva, P. S. Dester, H. Siljak, N. Marchetti, P. H. Nardelli, and R. A. de Souza, "Semantic-functional communications for multiuser

- event transmissions via random maps,” *arXiv preprint arXiv:2204.03223*, 2022.
- [34] X. Chen, T. Tan, G. Cao, and T. F. L. Porta, “Context-aware and energy-aware video streaming on smartphones,” *IEEE Transactions on Mobile Computing*, vol. 21, no. 3, pp. 862–877, Aug. 2020.
- [35] Z. Liu, Q. Z. Sheng, X. Xu, D. Chu, and W. E. Zhang, “Context-aware and adaptive qos prediction for mobile edge computing services,” *IEEE Transactions on Services Computing*, vol. 15, no. 1, pp. 400–413, Sept. 2019.
- [36] G. J. Stamatakis, N. Pappas, A. Fragkiadakis, and A. Traganitis, “Semantics-Aware Active Fault Detection in IoT,” in *Proc. of the 20th International Symposium on Modeling and Optimization in Mobile, Ad hoc, and Wireless Networks (WiOpt)*. Torino, Italy: IEEE, Sept. 2022, pp. 161–168.
- [37] E. Uysal, O. Kaya, A. Ephremides, J. Gross, M. Codreanu, P. Popovski, M. Assaad, G. Liva, A. Munari, B. Soret *et al.*, “Semantic communications in networked systems: A data significance perspective,” *IEEE Network*, vol. 36, no. 4, pp. 233–240, Oct. 2022.
- [38] A. Maatouk, M. Assaad, and A. Ephremides, “The age of incorrect information: An enabler of semantics-empowered communication,” *IEEE Transactions on Wireless Communications*, May 2022.
- [39] C. A. López and O. I. Lombardi, “No communication without manipulation: A causal-deflationary view of information,” *Studies in History and Philosophy of Science Part A*, vol. 73, pp. 34–43, Jun. 2019.
- [40] R. B. Ash, *Information theory*. Courier Corporation, 2012.
- [41] A. Kolchinsky and D. H. Wolpert, “Semantic information, autonomous agency and non-equilibrium statistical physics,” *Interface focus*, vol. 8, no. 6, p. 20180041, Sep. 2018.
- [42] O. Hashash, C. Chaccour, and W. Saad, “Edge continual learning for dynamic digital twins over wireless networks,” in *Proc. of the 23rd International Workshop on Signal Processing Advances in Wireless Communication (SPAWC)*, Oulu, Finland, Jun. 2022, pp. 1–5.
- [43] R. M. Gray, *Entropy and information theory*. Springer Science & Business Media, 2011.
- [44] D. Huang, X. Tao, F. Gao, and J. Lu, “Deep learning-based image semantic coding for semantic communications,” in *Proc. of IEEE Global Communications Conference (GLOBECOM)*, Madrid, Spain, Dec. 2021, pp. 1–6.
- [45] O. Gune, B. Banerjee, S. Chaudhuri, and F. Cuzzolin, “Generalized zero-shot learning using generated proxy unseen samples and entropy separation,” in *Proc. of the 28th ACM International Conference on Multimedia*, Seattle, WA, Oct. 2020, pp. 4262–4270.
- [46] J. Lecoq, M. Oliver, J. H. Siegle, N. Orlova, P. Ledochowitsch, and C. Koch, “Removing independent noise in systems neuroscience data using deepinterpolation,” *Nature methods*, vol. 18, no. 11, pp. 1401–1408, Nov. 2021.
- [47] A. Zhu, J. G. Ibrahim, and M. I. Love, “Heavy-tailed prior distributions for sequence count data: removing the noise and preserving large differences,” *Bioinformatics*, vol. 35, no. 12, pp. 2084–2092, Jun. 2019.
- [48] M. Fresia, F. Perez-Cruz, H. V. Poor, and S. Verdu, “Joint source and channel coding,” *IEEE Signal Processing Magazine*, vol. 27, no. 6, pp. 104–113, Oct. 2010.
- [49] N. Farsad, M. Rao, and A. Goldsmith, “Deep learning for joint source-channel coding of text,” in *Proc. of IEEE international conference on acoustics, speech and signal processing (ICASSP)*, Calgary, Canada, Apr. 2018, pp. 2326–2330.
- [50] E. Boursoulatzé, D. B. Kurka, and D. Gündüz, “Deep joint source-channel coding for wireless image transmission,” *IEEE Transactions on Cognitive Communications and Networking*, vol. 5, no. 3, pp. 567–579, May 2019.
- [51] A. Achille, G. Paolini, G. Mbeng, and S. Soatto, “The information complexity of learning tasks, their structure and their distance,” *Information and Inference: A Journal of the IMA*, vol. 10, no. 1, pp. 51–72, Mar. 2021.
- [52] J. Zilly, A. Achille, A. Censi, and E. Frazzoli, “On plasticity, invariance, and mutually frozen weights in sequential task learning,” in *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds., vol. 34. Curran Associates, Inc., 2021, pp. 12386–12399. [Online]. Available: <https://proceedings.neurips.cc/paper/2021/file/6738fc33dd0b3906cd3626397cd247a7-Paper.pdf>
- [53] H. Xie, Z. Qin, G. Y. Li, and B.-H. Juang, “Deep learning enabled semantic communication systems,” *IEEE Transactions on Signal Processing*, vol. 69, pp. 2663–2675, Apr. 2021.
- [54] Q. Zhou, R. Li, Z. Zhao, C. Peng, and H. Zhang, “Semantic communication with adaptive universal transformer,” *IEEE Wireless Communications Letters*, vol. 11, no. 3, pp. 453–457, Dec. 2021.
- [55] H. Zhao, R. T. Des Combes, K. Zhang, and G. Gordon, “On learning invariant representations for domain adaptation,” in *Proc. of the 36th International Conference on Machine Learning*, Long Beach, CA, Jun. 2019, pp. 7523–7532.
- [56] A. T. Nguyen, T. Tran, Y. Gal, and A. G. Baydin, “Domain invariant representation learning with domain density transformations,” in *Proc. of the 35th Conference on Neural Information Processing Systems*, vol. 34, Virtual, Dec. 2021, pp. 5264–5275.
- [57] Y. Lu and J. Lu, “A universal approximation theorem of deep neural networks for expressing probability distributions,” *Advances in neural information processing systems*, vol. 33, pp. 3094–3105, 2020.
- [58] Y. Yang, C. Guo, F. Liu, C. Liu, L. Sun, Q. Sun, and J. Chen, “Semantic communications with AI tasks,” *arXiv preprint arXiv:2109.14170*, 2021.
- [59] Z. Weng, Z. Qin, and G. Y. Li, “Semantic communications for speech signals,” in *IEEE International Conference on Communications*, Montréal, Canada, Jun. 2021, pp. 1–6.
- [60] H. Xie, Z. Qin, and G. Y. Li, “Task-oriented multi-user semantic communications for vqa,” *IEEE Wireless Communications Letters*, vol. 11, no. 3, pp. 553–557, Jul. 2021.
- [61] K. Lu, Q. Zhou, R. Li, Z. Zhao, X. Chen, J. Wu, and H. Zhang, “Rethinking modern communication from semantic coding to semantic communication,” *IEEE Wireless Communications*, May 2022.
- [62] J. Dai, S. Wang, K. Tan, Z. Si, X. Qin, K. Niu, and P. Zhang, “Nonlinear transform source-channel coding for semantic communications,” *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 8, pp. 2300–2316, Jun. 2022.
- [63] Y. Wang, M. Chen, T. Luo, W. Saad, D. Niyato, H. V. Poor, and S. Cui, “Performance optimization for semantic communications: An attention-based reinforcement learning approach,” *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 9, pp. 2598–2613, Jul. 2022.
- [64] F. Zhou, Y. Li, X. Zhang, Q. Wu, X. Lei, and R. Q. Hu, “Cognitive semantic communication systems driven by knowledge graph,” *arXiv preprint arXiv:2202.11958*, 2022.
- [65] S. Jiang, Y. Liu, Y. Zhang, P. Luo, K. Cao, J. Xiong, H. Zhao, and J. Wei, “Reliable semantic communication system enabled by knowledge graph,” *Entropy*, vol. 24, no. 6, p. 846, Jun. 2022.
- [66] Y. Wang, M. Chen, W. Saad, T. Luo, S. Cui, and H. V. Poor, “Performance optimization for semantic communications: An attention-based learning approach,” in *Proc. of IEEE Global Communications Conference (GLOBECOM)*, Madrid, Spain, Dec. 2021, pp. 1–6.
- [67] W. Babonnaud, “A topos-based approach to building language ontologies,” in *International Conference on Formal Grammar*. Springer, 2019, pp. 18–34.
- [68] J.-C. Belfiore and D. Bennequin, “Topos and stacks of deep neural networks,” *arXiv preprint arXiv:2106.14587*, 2021.
- [69] L. Heng and B. McColl, *Mathematics for Future Computing and Communications*. Cambridge University Press, 2021.
- [70] F. W. Lawvere and S. H. Schanuel, *Conceptual mathematics: a first introduction to categories*. Cambridge University Press, 2009.
- [71] B. Schölkopf and J. von Kügelgen, “From statistical to causal learning,” *arXiv preprint arXiv:2204.00607*, 2022.
- [72] B. Schölkopf, F. Locatello, S. Bauer, N. R. Ke, N. Kalchbrenner, A. Goyal, and Y. Bengio, “Toward causal representation learning,” *Proceedings of the IEEE*, vol. 109, no. 5, pp. 612–634, May 2021.
- [73] J. Pearl, *Causality*. Cambridge university press, 2009.
- [74] J. Kaddour, A. Lynch, Q. Liu, M. J. Kusner, and R. Silva, “Causal machine learning: A survey and open problems,” *arXiv preprint arXiv:2206.15475*, 2022.
- [75] C.-Y. Chuang, J. Robinson, Y.-C. Lin, A. Torralba, and S. Jegelka, “Debiased contrastive learning,” *Proc. of Advances in neural information processing systems*, vol. 33, pp. 8765–8775, Dec. 2020.
- [76] J. Mitrovic, B. McWilliams, J. Walker, L. Buesing, and C. Blundell, “Representation learning via invariant causal mechanisms,” *arXiv preprint arXiv:2010.07922*, 2020.
- [77] V. Veitch, A. D’Amour, S. Yadlowsky, and J. Eisenstein, “Counterfactual invariance to spurious correlations in text classification,” *Proc. of Advances in Neural Information Processing Systems*, vol. 34, pp. 16196–16208, Dec. 2021.
- [78] P. Smolensky, “Connectionist ai, symbolic ai, and the brain,” *Artificial Intelligence Review*, vol. 1, no. 2, pp. 95–109, 1987.
- [79] C. K. Thomas and W. Saad, “Neuro-symbolic causal reasoning meets signaling game for emergent semantic communications,” *arXiv preprint arXiv:2210.12040*, 2022.
- [80] E. Bengio, M. Jain, M. Korablyov, D. Precup, and Y. Bengio, “Flow network based generative models for non-iterative diverse candidate generation,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 27381–27394, Dec. 2021.
- [81] J. Liu, S. Shao, W. Zhang, and H. V. Poor, “An indirect rate-distortion characterization for semantic sources: General model and the case of gaussian observation,” *IEEE Transactions on Communications*, vol. 70, no. 9, pp. 5946–5959, Jul. 2022.
- [82] L. Gavrilovska, V. Rakovic, and D. Denkovski, “From cloud ran to open ran,” *Wireless Personal Communications*, vol. 113, no. 3, pp. 1523–1539, 2020.
- [83] M. Chen, U. Challita, W. Saad, C. Yin, and M. Debbah, “Artificial neural networks-based machine learning for wireless networks: A tutorial,” *IEEE Communications Surveys & Tutorials*, 2019.
- [84] Z. Han, D. Niyato, W. Saad, T. Başar, and A. Hjørungnes, *Game theory in wireless and communication networks: theory, models, and applications*. Cambridge university press, 2012.
- [85] N. S. Kovach, A. S. Gibson, and G. B. Lamont, “Hypergame theory: a model for conflict, misperception, and deception,” *Game Theory*, vol. 2015, 2015.
- [86] S. Niknam, A. Roy, H. S. Dhillon, S. Singh, R. Banerji, J. H. Reed, N. Saxena, and S. Yoon, “Intelligent o-ran for beyond 5g and 6g wireless networks,” *arXiv preprint arXiv:2005.08374*, 2020.
- [87] Intelligent Automation Guide Series. Intelligent optimization How intelligent RAN automation will re-energize the SON market. Stockholm, Sweden. [Online]. Available: <https://www.ericsson.com/49c540/assets/local/core-network/doc/intelligent-optimization-guide.pdf>