# Large Wireless Model (LWM): A Foundation Model for Wireless Channels

Sadjad Alikhani, Gouranga Charan, and Ahmed Alkhateeb

Wireless Intelligence Lab, Arizona State University, USA

Emails: {alikhani, gcharan, alkhateeb}@asu.edu

*Abstract*—This paper presents the Large Wireless Model (LWM)—the world's first foundation model for wireless channels. Designed as a task-agnostic model, LWM generates universal, rich, contextualized channel embeddings (features) that potentially enhance performance across a wide range of downstream tasks in wireless communication and sensing systems. Towards this objective, LWM, which has a transformer-based architecture, was pre-trained in a self-supervised manner on large-scale wireless channel datasets. Our results show consistent improvements in classification and regression tasks when using the LWM embeddings compared to raw channel representations, especially in scenarios with high-complexity machine learning tasks and limited training datasets. This LWM's ability to learn from large-scale wireless data opens a promising direction for intelligent systems that can efficiently adapt to diverse tasks with limited data, paving the way for addressing key challenges in wireless communication and sensing systems.

## I. INTRODUCTION

Current and future wireless communication and sensing systems feature important trends that promise substantial performance gains [2]–[4]. For example, these systems are rapidly relying on the use of large antenna arrays, the operation over high frequency bands in mid-band, millimeter wave (mmWave), and sub-terahertz, the support of massive number of communicating and sensing devices of various quality of service requirements, and the densification of network infrastructure nodes. Further, these wireless communication and sensing systems increasingly interact with each other, from coordination and integration to assisting each other. Achieving the high potential of these new trends, however, requires overcoming critical challenges in high-dimensional signal processing, complex optimization problems, massive wireless overhead requirements, and intricate network management among others. All that motivates the development of novel approaches for the modeling, optimizing, and operation of next-generation wireless communication and sensing systems.

Traditional modeling techniques, such as statistical models and optimization-based approaches, struggle to address these challenges effectively. These methods often rely on simplified models or scenario-specific features, failing to generalize across the diverse and dynamic environments of future wireless communication and sensing networks. For instance, they may not capture complex interference patterns in dense small-cell networks or scale poorly to high-dimensional MIMO systems. Deep learning has emerged as a promising alternative, offering data-driven (model-based or model-free) solutions for optimizing network performance, resource allocation, and signal processing [5]–[9]. However, deep learning approaches also face significant limitations. First, they typically **require large labeled datasets**, which are often scarce in wireless networks and are typically expensive and hard to collect. Second, traditional deep learning models like convolutional neural networks (CNNs) and recurrent neural networks (RNNs) struggle with specific aspects of wireless communication and sensing tasks. CNNs may not capture temporal dependencies efficiently [10], [11], while RNNs often struggle with long-term dependencies and real-time computational efficiency [1], [12]. These limitations underscore the need for a more robust and adaptable framework for leveraging and deploying deep learning in wireless communication and sensing networks.

To address these challenges, we propose Large Wireless Model (LWM), a foundation model specifically designed for wireless communication and sensing channels. LWM introduces a task-agnostic framework with pre-training on large-scale synthetic data. As a task-agnostic model, **LWM serves as a universal feature extractor for multiple downstream tasks, facilitating complex problem-solving with limited labeled data**. It leverages transformer models with multi-head attention mechanisms to capture complex spatial and temporal relationships in wireless channel data. Inspired by advancements in natural language processing (NLP) [1], [13]–[16], audio processing [17]–[19], and computer vision [20]–[22], LWM learns rich, context-aware embeddings that can be utilized for various downstream wireless tasks, such as channel estimation, beamforming, and interference management. LWM is *pre-trained* on extensive wireless channel datasets, covering a wide range of wireless scenarios. This approach enables the model to capture fundamental properties of wireless propagation and network dynamics, which can be transferred to real-world scenarios, even with limited task-specific data. Through these innovations, LWM addresses the key challenges of limited labeled data, complex spatial-temporal dependencies, and the need for generalization across diverse wireless environments. The key contributions can be summarized as follows:

- We introduce the world's first foundation model for wireless channel embeddings, capable of extracting universal, rich, and context-aware features from complex wireless channels and in diverse environments.
- We demonstrate the LWM's effectiveness across multiple downstream tasks, showcasing its ability to generalize to
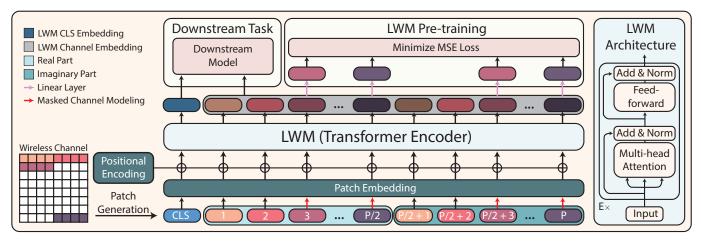
Figure 1: This figure depicts the offline pre-training and online embedding generation process for LWM. The channel is divided into fixed-size patches, which are linearly embedded and combined with positional encodings before being passed through a Transformer encoder. During self-supervised pre-training, some embeddings are masked, and LWM leverages self-attention to extract deep features, allowing the decoder to reconstruct the masked values. For downstream tasks, the generated LWM embeddings enhance performance. The right block shows the LWM architecture, inspired by [1].

various wireless scenarios with limited task-specific data.

- We provide a comprehensive analysis of the LWM's performance compared to conventional approaches that use raw wireless channels, highlighting its advantages in feature extraction and generalization.

This work introduces a new framework for leveraging and deploying deep learning in wireless communication and sensing systems by leveraging the power of foundation models to address key modeling, design, and deployment challenges in wireless systems. The pre-trained LWM model, scripts, datasets, demo, and instructions are available on the Wireless Intelligence Lab's Hugging Face page [1], allowing researchers to incorporate them into their projects.

## II. PRIOR WORK

Foundation models have transformed artificial intelligence by introducing a paradigm of large-scale pre-training followed by task-specific fine-tuning [14]–[20]. These models, exemplified by Bidirectional Encoder Representations from Transformers (BERT) [13] in NLP and Wav2Vec 2.0 [17] in audio processing, leverage transformer architectures with multi-head attention mechanisms to capture complex relationships in data. The pre-training phase typically involves massive datasets and novel learning objectives: BERT uses self-supervised methods, namely masked language modeling and next-sentence prediction, while Wav2Vec 2.0 employs contrastive learning and masked prediction tasks for audio understanding. This process allows the models to learn rich, contextual representations of their input domains. The key to their success lies in the attention mechanism, which enables dynamic focus on relevant parts of the input, and the scale of pre-training, which allows the model to capture a wide range of patterns and relationships. The resulting pre-trained model serves as a powerful feature extractor, capable of being fine-tuned on various downstream

tasks with limited task-specific data. This transfer learning capability, combined with the model's ability to capture long-range dependencies and generalize across scenarios, has led to state-of-the-art performance across numerous applications.

Unlike traditional models, which process sequences sequentially (like RNNs) or in a restricted local context (like CNNs), Transformers look at all parts of a sequence simultaneously. This parallel processing allows each element to relate to every other element, capturing dependencies across the entire input sequence. For example, in a sentence, words like "bank" and "river" might appear far apart, but self-attention lets the model understand their association when the context implies a natural setting, rather than financial. By calculating attention scores between each pair of words, Transformers allow each word to dynamically "attend" to others, helping the model to understand nuanced meanings and relationships.

In language models like BERT, each word token in a sentence is converted into a dense vector representation (embedding) and is associated with a query, key, and value. The attention mechanism then calculates a weighted average of these values for each word, where the weights are derived from the similarity between the queries and keys. This process allows the model to focus on important words and phrases depending on context. For instance, in the sentence, "The animal didn't cross the street because it was too tired," the model can determine that "it" refers to "the animal" rather than "the street." This disambiguation is achieved because self-attention scores between "it" and "the animal" would be higher than between "it" and "the street."

Another example can be seen in GPT models [16], where the attention mechanism not only understands contextual nuances but also generates coherent text. Given a prompt, the model iteratively predicts the next word by attending to all previous words, ensuring that each generated word maintains context. For instance, when prompted with "Once upon a time in a small village," GPT can generate a follow-up that maintains

the narrative theme by giving more weight to probable story elements over unrelated concepts. This process enables the model to construct coherent, contextually relevant text, showcasing the power of the self-attention mechanism to understand complex dependencies and generate highly fluent and context-aware language output.

Moreover, in vision models, Transformers leverage the same self-attention mechanism to capture spatial relationships across pixels or image patches, providing a powerful alternative to convolutional neural networks (CNNs). For example, in Vision Transformers (ViTs) [20], an image is divided into small, fixed-size patches, each of which is flattened and linearly embedded, similar to tokens in a sentence. These embeddings are then processed in parallel, allowing the Transformer to learn relationships between distant regions of an image. This capability is particularly useful for tasks requiring an understanding of global context, such as recognizing objects within complex backgrounds or interpreting spatial patterns across an entire scene. For instance, in a complex scene containing a forest with animals partially obscured by trees, CNNs might struggle to understand the spatial relationship between scattered animal parts due to their limited receptive fields. However, a ViT can relate these distant regions by directly computing attention scores between patches, allowing it to "see" the entire animal even if parts of it are distant in pixel space. This method enables ViTs to recognize high-level patterns and achieve state-of-the-art performance in tasks like object detection, image segmentation, and visual question answering, where a global understanding of the image context is essential.

The success of foundation models in other domains presents compelling opportunities for **wireless communications and sensing**, particularly in addressing challenges related to complex spatial-temporal dependencies and limited labeled data. However, adapting this paradigm to wireless systems requires overcoming several hurdles. Unlike text or images, wireless signals have unique characteristics such as complex-valued data, rapid temporal variations, and domain-specific noise patterns. The lack of large-scale, diverse datasets in wireless communications comparable to those in NLP or computer vision poses another challenge. Despite these obstacles, a foundation model for wireless communications and sensing could potentially provide a **universal feature extractor** for several tasks. By pre-training on extensive datasets (either synthetic datasets generated through advanced ray-tracing simulations or real-world data), such a model could capture fundamental properties of wireless propagation and network dynamics. This approach could enable robust performance across diverse wireless environments, even with limited task-specific data, addressing key challenges in the field and potentially revolutionizing the field.

## III. LWM

Building upon the foundation model paradigm discussed in Section II, LWM applies these principles to the domain of wireless communications (and sensing). LWM is designed as a task-agnostic model for generalized feature extraction in wireless channels, addressing the unique challenges of this field. LWM adopts a self-supervised Transformer architecture with multi-head attention mechanisms and is pre-trained on a large dataset of wireless channels. LWM processes input channels in patches, enabling it to serve as a universal feature extractor for various wireless communication and sensing tasks. This approach allows LWM to capture complex patterns and provide rich, contextual representations of wireless environments, potentially improving performance across diverse scenarios even with limited task-specific data.

LWM is built upon several core design principles:

**Patch-Based Processing:** Wireless channels are segmented into patches, enabling LWM to capture both local and global patterns efficiently. This patch-based structure allows for spatial and spectral dependencies to be encoded in a way that mimics human perception of relevant wireless features, enhancing LWM's utility as a universal feature extractor.

**Self-Supervised Pre-Training:** Unlike traditional supervised models that require labeled data for each task, LWM is trained on a large dataset of unlabeled wireless channels using self-supervised techniques. By leveraging masked channel modeling and an attention mechanism, LWM learns to capture complex structural relationships within the data without relying on labeled datasets.

**Multi-Head Attention:** The attention mechanism in LWM allows it to selectively focus on relevant patterns in wireless channels, dynamically assigning importance to different parts of the input. This is particularly beneficial for wireless applications, where meaningful features can vary spatially and spectrally across environments.

**Task Flexibility and Transferability:** The representations produced by LWM are highly versatile, making it possible to apply the model to numerous downstream tasks with minimal or no fine-tuning. As a result, LWM can generalize across different scenarios and geographic regions, achieving robust performance even when task-specific data is sparse or highly variable.

In the following sections, we detail each component of LWM's pipeline, from data preprocessing and embedding to the architecture of the Transformer encoder blocks. These sections will demonstrate how LWM processes raw wireless channels into enriched feature representations and embeddings, as illustrated in Fig. 1, and provide insight into how this model bridges the gap between traditional wireless communication techniques and advanced, general-purpose AI architectures.

## IV. DATA PREPROCESSING: MASKING AND EMBEDDING

Given the input wireless channels, we first preprocess them to align with the Transformer input format and facilitate our self-supervised pre-training method. The steps for this are outlined below.

### A. Patch Generation

Unlike RNNs, Transformers require the entire input simultaneously. To achieve this, we feed each channel in a patch-based

format [20]. Each channel matrix $\mathbf{H} \in \mathbb{C}^{M \times N}$ is split into $P$ patches by first separating the real and imaginary components, flattening them, and then dividing each part into patches. The process is as follows:

1) **Separate Real and Imaginary Components:** We start by separating $\mathbf{H}$ into its real and imaginary components

$$\mathbf{H}_{\text{real}} = \Re(\mathbf{H}), \tag{1a}$$
$$\mathbf{H}_{\text{imag}} = \Im(\mathbf{H}), \tag{1b}$$

where $\mathbf{H}_{\text{real}}, \mathbf{H}_{\text{imag}} \in \mathbb{R}^{M \times N}$.

2) **Flatten Each Component:** Flatten both $\mathbf{H}_{\text{real}}$ and $\mathbf{H}_{\text{imag}}$ into vectors using the vectorization operator, as follows

$$\mathbf{h}_{\text{real}} = \text{vec}\left(\mathbf{H}_{\text{real}}^{\top}\right), \tag{2a}$$
$$\mathbf{h}_{\text{imag}} = \text{vec}\left(\mathbf{H}_{\text{imag}}^{\top}\right), \tag{2b}$$

resulting in $\mathbf{h}_{\text{real}}, \mathbf{h}_{\text{imag}} \in \mathbb{R}^{MN}$.

3) **Divide into Patches:** Divide each flattened component into $P/2$ patches. Let each patch have a length $L = 2MN/P$. We generate patches as follows

$$\mathbf{p}_i = \mathbf{h}_{\text{real}}[(i-1)L + 1 : iL], \quad i \in [P/2], \tag{3a}$$
$$\mathbf{p}_{i+P/2} = \mathbf{h}_{\text{imag}}[(i-1)L + 1 : iL], \quad i \in [P/2], \tag{3b}$$

where each patch $\mathbf{p}_i \in \mathbb{R}^L$, with $i \in [P] = \{1, 2, \ldots, P\}$.

This patch-based approach accelerates computations, enables the model to learn both inter- and intra-patch relationships, mimics convolutional layers with self-attention, and increases design flexibility. For instance, with $M = 32$, $N = 32$, $P = 128$, and $L = 16$, this configuration generates 128 patches of length 16. Smaller patches capture fine-grained details but require higher computational cost due to the increased number of patches, and they may limit the model's ability to learn broader structural dependencies [23]. Larger patches emphasize long-range dependencies and structural patterns, reducing computational requirements but potentially overlooking finer details. Ultimately, selecting the optimal patch size depends on the task's need for local versus global information and the available computational resources, rather than any specific threshold.

For LWM, which is designed to serve as a universal feature extractor, achieving this balance is essential. By selecting a patch size that allows both detail and structure to be preserved, LWM remains adaptable, capturing features that are sufficiently detailed to represent nuanced variations while broad enough to generalize across tasks and datasets. A useful approach for assessing an optimal patch size is to evaluate the performance of the masking strategy described next.

### B. Masked Channel Modeling

For LWM to be task-agnostic, we pre-train it in a self-supervised manner. Self-supervised learning enables LWM to capture the intrinsic structure of the data, allowing it to serve as a universal feature extractor adaptable to various downstream tasks, without reliance on labeled data. To enable this self-supervised pre-training, we propose a technique called Masked Channel Modeling (MCM).

In MCM, we mask 15% of the real part patches. This percentage is chosen to balance the model's access to context with the challenge of reconstructing missing patches. Masking too many patches could lead to excessive information loss, making reconstruction difficult, while masking too few would provide limited incentive for the model to learn complex dependencies. Within the selected 15% of patches, specific sub-percentages are assigned:

- 80% **are fully masked** with a uniform vector $\mathbf{m} = [m, m, \ldots, m]^{\top} \in \mathbb{R}^L$, preventing the model from accessing information in those locations. This requires LWM to leverage surrounding patches to predict the masked values, pushing it to learn the spatial dependencies within the channel.
- 10% **are replaced with random vectors** sampled from a distribution (e.g., $\mathcal{N}(0, \sigma^2)$), adding noise and encouraging LWM to differentiate genuine channel structures from anomalies.
- 10% **are left unchanged**, providing partial ground truth to stabilize the model's predictions and helping it recognize real patterns among masked and altered patches.

This masking strategy is carefully applied to both the real and imaginary parts to prevent information leakage between them. Specifically, the imaginary patches selected for masking are the exact counterparts of the randomly selected real patches. This ensures that if a real part patch is masked, the corresponding imaginary part patch is also masked, preventing any indirect inference. If only the real part of a patch $\mathbf{p}_i$ were masked while leaving its imaginary part $\mathbf{p}_{i+P/2}$ unmasked, the model could leverage the unmasked imaginary patch $\mathbf{p}_{i+P/2}$ to predict $\mathbf{p}_i$. By masking both components of each selected patch, LWM learns each component's structure independently, resulting in robust feature extraction and reducing the chance of unintended data leakage.

The selected 15% of patches are masked before going through the input embedding, positional encoding, and finally the LWM pre-training stages. The corresponding high-dimensional embeddings of masked patches at the output of LWM then pass through a simple linear layer that maps each embedding back to the original patch size. The goal is to minimize the Mean Squared Error (MSE) between the reconstructed masked patches and their original values, expressed as

$$\mathcal{L}_{\text{MCM}} = \frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \left\| \mathbf{W}_i^{\text{dec}} \mathbf{e}_i^{\text{LWM}} - \mathbf{p}_i \right\|^2, \tag{4}$$

where $\mathcal{M}$ represents the set of all selected (masked) patches, $\mathbf{e}_i^{\text{LWM}} \in \mathbb{R}^D$ is the high-dimensional embedding of the $i$-th masked patch at LWM's output, $\mathbf{W}_i^{\text{dec}} \in \mathbb{R}^{L \times D}$ is the weight matrix of the linear layer used to map $\mathbf{e}_i$ back to the original patch size, and $\mathbf{p}_i$ is the original value of the $i$-th patch.

This approach allows LWM to develop highly refined embeddings that can be decoded accurately using a simple linear layer, underscoring the richness and expressiveness of the learned representations. Because the Transformer encoder does

not know which patches will be masked or replaced by random patches, it is forced to maintain a contextual representation for every input patch, ensuring the embeddings are robust and contextually aware. Additionally, as random replacement only occurs for a small fraction of all patches (10% of the 15% masked patches), the model's ability to capture spatial and structural dependencies remains unaffected [13]. The high performance of these linear layers in reconstructing masked patches highlights the quality of LWM's output embeddings, which effectively capture complex spatial relationships within the channel, enabling efficient decoding and ensuring robust performance across a variety of downstream tasks.

### C. CLS Patch

We prepend an additional patch, known as the CLS (classification) patch [13], to the sequence of channel patches, increasing the sequence length to $P+1$. The CLS patch is initialized as a learnable vector, denoted $\mathbf{p}_{\text{CLS}} = [c, c, \ldots, c]^\top \in \mathbb{R}^L$, where $c$ is set as a random value. Through its interactions across the Transformer layers, it aggregates and summarizes information from all other patches in the sequence.

This interaction enables the CLS patch to capture a comprehensive view of the entire input by attending to each patch in each layer, accumulating information from both local details and broader structures. As the sequence passes through multiple Transformer layers, the CLS patch's representation is refined, with each layer integrating lower-level details (e.g., spatial or temporal variations) and higher-level features (e.g., overall channel quality or dominant paths). The resulting representation, $\mathbf{e}_{\text{CLS}}^{\text{LWM}}$, serves as a compact, high-level summary of the input sequence

$$\mathbf{e}_{\text{CLS}}^{\text{LWM}} = f_{\text{LWM}}(\mathbf{e}_{\text{CLS}}, \{\mathbf{e}_i\}_{i=1}^P), \quad (5)$$

where $f_{\text{LWM}}$ represents the transformations and interactions within LWM. Through this mechanism, the CLS patch identifies the patches that contribute most significantly to the structure and dependencies within the channel, highlighting the segments most critical in defining the channel's overall characteristics.

Additionally, the compact size of the CLS patch makes it an efficient, low-dimensional encoded representation of the channels. This compactness is advantageous, as it serves as an expressive summary without requiring further model training, given that LWM has already been pre-trained in a task-agnostic manner to capture these global features.

The **attention scores** assigned between the CLS patch and each patch $\mathbf{p}_i$ reveal the relative importance of each patch, making the CLS patch a valuable interpretability tool across various wireless tasks. For example, in **channel estimation**, the CLS patch can highlight critical segments of the channel, such as those representing major multipath components, essential for accurate channel estimation. In **classification tasks** such as line-of-sight (LoS) and non-line-of-sight (NLoS) classification, the CLS patch effectively summarizes the input, capturing the distinguishing features needed to identify LoS versus NLoS conditions. By attending to patches with high

relevance to signal paths, reflection, and scattering, the CLS patch representation $\mathbf{e}_{\text{CLS}}^{\text{LWM}}$ captures a robust global context that is ideal for classification tasks, where understanding the overall channel state is crucial. Additionally, in **resource allocation**, the CLS patch's attention scores $\{\alpha_{\text{CLS},i}\}_{i=1}^P$ (where $\alpha_{\text{CLS},i}$ is the attention weight of the $i$-th patch with respect to the CLS patch) help identify channel segments that demand more resources. This optimization supports spectrum and power allocation by prioritizing the most influential patches. In **beamforming** and **beam selection**, the CLS patch can focus on patches that carry directional cues, assisting in adaptive beam selection by highlighting segments most indicative of optimal beam configurations. In **semantic communications**, the CLS patch enables efficient encoding by capturing patches containing the most contextually relevant information, enhancing communication quality while reducing redundancy. These examples demonstrate the versatility of the **CLS patch** across wireless applications, where it serves as an adaptable **focal point** for understanding and prioritizing channel features that directly impact performance [24].

### D. Input Embedding

After masking and prepending the CLS patch to the channel patch sequence, each patch is projected into an embedding space with dimension $D$ using a linear layer, effectively mapping each flattened patch to a $D$-dimensional vector. Given a set of patches $\{\mathbf{p}_{\text{CLS}}, \mathbf{p}_1^m, \mathbf{p}_2^m, \ldots, \mathbf{p}_P^m\}$, where $\mathbf{p}_i^m \in \mathbb{R}^L$ represents a masked patch (with only 15% of patches masked), the linear layer performs the following transformation for each patch $\mathbf{p}_i^m$

$$\mathbf{e}_i^{\text{emb}} = \mathbf{W}_i^{\text{emb}} \mathbf{p}_i^m + \mathbf{b}_i \in \mathbb{R}^D, \quad i \in \{\text{CLS}\} \cup [P], \quad (6)$$

where $\mathbf{W}_i^{\text{emb}} \in \mathbb{R}^{D \times L}$ is the weight matrix, and $\mathbf{b}_i \in \mathbb{R}^D$ is the bias vector for the $i$-th patch. This transformation produces the initial patch embeddings $\mathbf{E}^{\text{emb}} = [\mathbf{e}_{\text{CLS}}^{\text{emb}}, \mathbf{e}_1^{\text{emb}}, \mathbf{e}_2^{\text{emb}}, \ldots, \mathbf{e}_P^{\text{emb}}]^\top \in \mathbb{R}^{(P+1) \times D}$, allowing the Transformer to process all patches in a common high-dimensional feature space where relationships can be effectively captured.

Embedding patches into a higher-dimensional space before feeding them into the Transformer is essential for capturing complex relationships, especially compared to models like autoencoders, which often reduce dimensionality. A higher-dimensional embedding provides a richer, more expressive feature space, enabling each patch to retain fine-grained information about the channel data. Directly using the original patch size as an embedding would limit the **model's capacity** to capture meaningful relationships. The choice of a higher embedding dimension $D$ enables the Transformer to establish detailed contextual relationships, similar to how embeddings in text-based models capture the semantic relationships between words. In language models, embeddings allow each token (word or subword) to represent not only its isolated meaning but its meaning contextualized by surrounding words. This context-awareness is essential for handling nuances like polysemy.

Likewise, in wireless channels, high-dimensional embeddings allow each patch to capture its information with respect to the entire channel, representing dependencies between different parts of the channel and capturing structural nuances. In wireless channels, fine-grained details such as spatial correlations, multipath effects, and scattering are crucial for accurate representation. By embedding patches into a higher dimension, we ensure sufficient capacity to capture these intricate dependencies. Lower-dimensional representations, as in autoencoders, often sacrifice such details for compression, which can limit the Transformer's capacity to understand the channel's complex structure.

Furthermore, **positional encodings** are added to these embeddings to provide the Transformer with information about the order of patches, as it inherently lacks sequence awareness. To achieve this, we define a positional encoding patch $\mathbf{p}_i^{\text{pos}} \in \mathbb{R}^L$ for each patch $i$ as follows:

- For the CLS token, the positional encoding patch is a uniform vector of zeros.
- For each subsequent patch $i \in [P]$, the positional encoding patch is a uniform vector filled with the value $i$, providing an incremental encoding. Mathematically, this is expressed as

$$\mathbf{p}_i^{\text{pos}} = i \cdot \mathbf{1}_L, \quad i \in [P], \tag{7}$$

where $\mathbf{1}_L$ is a vector of ones in $\mathbb{R}^L$.

Each positional encoding patch $\mathbf{p}_i^{\text{pos}}$ is then mapped into the embedding space using a learned embedding matrix $\mathbf{W}_i^{\text{pos}} \in \mathbb{R}^{D \times L}$ and a bias vector $\mathbf{b}_i^{\text{pos}} \in \mathbb{R}^D$, resulting in positional encodings

$$\mathbf{e}_i^{\text{pos}} = \mathbf{W}_i^{\text{pos}} \mathbf{p}_i^{\text{pos}} + \mathbf{b}_i^{\text{pos}}, \quad i \in \{\text{CLS}\} \cup [P]. \tag{8}$$

These position embeddings are then added to the corresponding patch embeddings $\mathbf{e}_i$, yielding the position-encoded input embeddings

$$\mathbf{e}_i^{\text{input}} = \mathbf{e}_i^{\text{emb}} + \mathbf{e}_i^{\text{pos}}, \quad i \in \{\text{CLS}\} \cup [P]. \tag{9}$$

This addition forms the final set of position-encoded embeddings $\mathbf{E}^{\text{input}} = [\mathbf{e}_{\text{CLS}}^{\text{input}}, \mathbf{e}_1^{\text{input}}, \mathbf{e}_2^{\text{input}}, \ldots, \mathbf{e}_P^{\text{input}}]^\top \in \mathbb{R}^{(P+1) \times D}$, effectively incorporating ordered context into each patch embedding for the Transformer model.

## V. Model Architecture

We adopt the Transformer encoder architecture from [1], with modifications tailored to the context of wireless channels. The model processes input embeddings through a sequence of $E$ encoder blocks, progressively refining the embeddings to capture increasingly complex relationships and contextual patterns in the data. Across the $n$-th encoder block, where $n \in [E]$, the input embeddings $\mathbf{E}_n^{\text{input}} \in \mathbb{R}^{(P+1) \times D}$ evolve, with each layer enhancing feature richness. This sequence culminates in the final output embeddings, denoted as $\mathbf{E}^{\text{LWM}} \in \mathbb{R}^{(P+1) \times D}$, which constitute the LWM embeddings, encapsulating a refined representation suitable for downstream tasks.

Within each encoder block, the input embeddings go through the following components sequentially: **multi-head attention**, **layer normalization** with **residual connections**, a **feed-forward network**, and another layer normalization with residual connections. Each component refines the embeddings by capturing different aspects of the data structure. The output from one encoder block serves as the input to the next block, enabling the model to progressively build complex representations that capture the dependencies within the wireless channels. The details of these model components *at each encoder block* are outlined below.

### A. Self-Attention Mechanism

The self-attention mechanism allows each patch in the input sequence to assign contextually weighted importance to all other patches, achieved through a dot-product similarity measure. Given input embeddings $\mathbf{E}^{\text{input}} \in \mathbb{R}^{(P+1) \times D}$, where $P$ is the number of patches and $D$ is the embedding dimension, each row of $\mathbf{E}^{\text{input}}$ represents an embedding vector corresponding to a patch.

To compute self-attention, we first derive the *Query* ($\mathbf{Q}$), *Key* ($\mathbf{K}$), and *Value* ($\mathbf{V}$) matrices using linear transformations, each defined by learned weights

$$\mathbf{Q} = \mathbf{E}^{\text{input}} \mathbf{W}^Q, \quad \mathbf{K} = \mathbf{E}^{\text{input}} \mathbf{W}^K, \quad \mathbf{V} = \mathbf{E}^{\text{input}} \mathbf{W}^V, \tag{10}$$

where $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{D \times D'}$. Here, $D'$ is typically set to $D$ for single-head attention but may vary to $D_{\text{H}}$ for multi-head attention, as will be discussed next.

The core of self-attention begins by calculating the **scaled dot-product** of the query ($\mathbf{Q}$) and key ($\mathbf{K}$) matrices, capturing the similarity among patches in the input sequence. This similarity measure is computed as follows

$$\mathbf{S} = \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D'}}, \tag{11}$$

where $\mathbf{S}$ represents the scaled similarity scores, and the scaling factor $\sqrt{D'}$ prevents gradient issues like explosion or vanishing. These scaled similarity scores are then normalized by applying the softmax function

$$\mathbf{A} = \text{softmax}(\mathbf{S}), \tag{12}$$

where $\mathbf{A}$ is the attention weight matrix, which ensures that the weights assigned to each patch sum to 1 for each query, allowing the model to focus selectively on the most relevant patches. Finally, the attention weights are used to compute a weighted sum of the value ($\mathbf{V}$) matrix, producing the attention output

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{A}\mathbf{V}. \tag{13}$$

This process enables the model to dynamically adjust its focus based on the contextual relevance of each patch in relation to others, effectively capturing both local and global relationships across the input sequence.

The Query, Key, and Value matrices in the self-attention mechanism represent distinct, interrelated aspects of each input patch, working together to capture both local and global dependencies in the data. Each of these matrices has a specific purpose, ultimately contributing to how the model emphasizes certain relationships over others within the input sequence.

**Query (Q):** The Query matrix encodes the **intent** or **interest** of each patch in the sequence, representing the way each patch seeks information relevant to itself from other patches. Queries help determine what aspects of the surrounding data are most pertinent to the current patch, as each query vector in $\mathbf{Q}$ essentially acts as a **question** that asks how similar or relevant other patches are to it.

**Key (K):** The Key matrix complements the Query matrix by encoding **characteristics** of each patch that can be evaluated by queries from other patches. While Queries seek information, Keys provide the attributes or features by which each patch can be "queried." For any two patches, their query and key vectors are compared to assess their mutual relevance, which is captured through the dot product $\mathbf{Q}\mathbf{K}^\top$. High similarity scores indicate that certain patches carry information of high interest to one another.

**Value (V):** The Value matrix represents the actual information that each patch provides to the model, contributing to the weighted sum in the final self-attention calculation. Once the model identifies relevant patches (through the similarity scores from the dot product of Queries and Keys), the corresponding Value vectors are aggregated based on their importance, defined by the computed attention weights. Thus, Values are akin to the **content** that is passed along in the attention mechanism, shaping the contextualized representation of each patch by integrating relevant details from others.

To better understand the roles of Query, Key, and Value, consider the analogy of a conference setting. Imagine each participant has specific interests or questions (Queries) they want addressed, such as AI techniques or wireless advancements. Each participant also has unique expertise (Keys) they can offer. If someone's interest in AI aligns with another's data science expertise, a strong match occurs, enabling an exchange of knowledge (Values). In self-attention, this relationship is formalized by computing the similarity between each Query and all Keys. When a Query and Key pair strongly align, a higher attention weight is assigned, resulting in a weighted aggregation of the corresponding Value vectors. This effectively allows each patch to gather information most relevant to its context.

Together, Queries, Keys, and Values facilitate a process where each patch determines which other patches to **attend** to, assigns importance to each based on their contextual relevance, and then updates its representation based on a weighted sum of the relevant information. This arrangement allows the Transformer model to dynamically focus on relationships across patches, enabling it to capture both immediate and extended dependencies in the data with great flexibility and precision.
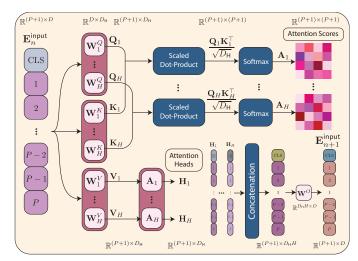


Figure 2: Multi-head Attention Mechanism

### B. Multi-Head Attention

Multi-head attention builds upon self-attention by enabling the model to learn multiple representations simultaneously, as depicted in Fig. 2. Instead of a single attention head, $H$ independent self-attention mechanisms are applied, each with distinct learned weight matrices

$$\text{head}_h = \text{Attention}(\mathbf{E}^{\text{input}}\mathbf{W}_h^Q, \mathbf{E}^{\text{input}}\mathbf{W}_h^K, \mathbf{E}^{\text{input}}\mathbf{W}_h^V), \quad (14)$$

where $\mathbf{W}_h^Q, \mathbf{W}_h^K, \mathbf{W}_h^V \in \mathbb{R}^{D \times D_\mathsf{H}}$ are the transformation matrices specific to each head, and $D_\mathsf{H} = \lfloor D/H \rfloor$. By processing across multiple heads, the model learns representations from various subspaces of the data, enhancing its ability to capture nuanced relationships. The outputs of all heads are concatenated and passed through a linear transformation

$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\text{head}_1, \ldots, \text{head}_H)\mathbf{W}^O, \quad (15)$$

where $\mathbf{W}^O \in \mathbb{R}^{D_\mathsf{H}H \times D}$. This combined output from multiple heads enables the model to focus on diverse aspects of the input sequence simultaneously, enriching the feature representation by capturing different patches and interactions in parallel.

In other words, each attention head can be thought of as a unique **lens** through which the model interprets the input data. In a multi-faceted problem space, these different lenses enable the model to pick up on varying levels of detail and relational structures simultaneously. To extend our analogy of participants in a conference: imagine that each participant not only has a set of questions (Queries) and expertise (Keys) but can also approach the conversation from multiple viewpoints or specializations (e.g., technical details, future trends, industry applications). One head might focus on broad concepts (like general channel structure or coarse spatial relationships), while another might pick up fine-grained, detailed aspects (such as precise timing or frequency patterns).

## C. Feed-Forward Network (FFN)

Each encoder layer has a feed-forward network (FFN) that processes individual patch embeddings independently. Formally, the FFN applies two linear transformations separated by a non-linear ReLU activation

$$\text{FFN}(\mathbf{e}_i) = \max(\mathbf{0}, \mathbf{e}_i \mathbf{W}_{i,1} + \mathbf{b}_{i,1})\mathbf{W}_{i,2} + \mathbf{b}_{i,2}, \qquad (16)$$

where $\mathbf{W}_{i,1} \in \mathbb{R}^{D \times D_{\text{FF}}}$ and $\mathbf{W}_{i,2} \in \mathbb{R}^{D_{\text{FF}} \times D}$. The intermediate layer size $D_{\text{FF}} = TD$ expands the embedding dimension by a factor of $T$, which empirically improves the model's expressiveness. The FFN enriches each embedding with more complex transformations, and because each embedding is processed separately, the FFN layer does not introduce dependencies across patches.

## D. Layer Normalization and Residual Connections

Layer normalization stabilizes the training by ensuring that each layer's output has zero mean and unit variance

$$\text{LayerNorm}(\mathbf{e}_i) = \frac{\mathbf{e}_i - \mu_i}{\sigma_i}, \qquad (17)$$

where $\mu_i$ and $\sigma_i$ denote the mean and standard deviation across the feature dimension of each patch embedding $\mathbf{e}$. This normalization is essential for convergence, especially in deeper networks, as it reduces covariate shifts.

Residual connections are added around each sub-layer (self-attention and FFN) to improve gradient flow

$$\text{Output}_i = \text{LayerNorm}(\mathbf{e}_i + \text{sub-layer}(\mathbf{e}_i)). \qquad (18)$$

The addition operation allows gradient information to flow more directly through the network, mitigating vanishing gradient problems and enabling efficient backpropagation.

## VI. PRE-TRAINING

Given the data preprocessing steps and the LWM model architecture described in the previous two subsections, we pre-train this model to be then leveraged in multiple downstream tasks. The following subsections outline the key components of LWM's pre-training process.

### A. Pre-Training Dataset

The LWM pre-training process utilizes a vast and diverse dataset of over 1 million wireless channels from 15 scenarios within the DeepMIMO dataset [25], with $80\%$ of the data dedicated to training and $20\%$ to validation. These scenarios include *O1*, *Boston5G*, *ASU Campus*, *New York*, *Los Angeles*, *Chicago*, *Houston*, *Phoenix*, *Philadelphia*, *Miami*, *Dallas*, *San Francisco*, *Austin*, *Columbus*, and *Seattle*. The first three scenarios are significantly larger, while the remaining twelve city scenarios average around 1500 effective users each. In wireless communications and sensing, relationships within the data often manifest as dependencies across various dimensions, requiring the model to capture both local patterns and broader spatial structures. Transformers are most effective when trained on large-scale and diverse datasets [20], as these enable them to generalize and learn complex relationships.

Training on this extensive dataset ensures that LWM captures nuanced wireless patterns, making it highly capable for a variety of downstream tasks.

### B. Wireless Channel Adaptation for Transformer Architecture

The wireless channels in the adopted dataset are structured as $(M, N) = (32, 32)$ matrices, where $N$ corresponds to the number of subcarriers and $M$ to the number of antennas. To prepare this data for Transformer processing, each channel matrix $\mathbf{H} \in \mathbb{C}^{32 \times 32}$ is split into $P = 128$ patches: 64 from the real part and 64 from the imaginary part, as detailed in section (IV). Each patch contains values from $L = 16$ consecutive subcarriers out of the 32 total subcarriers, spanning each antenna, resulting in patches of size $16 \times 1$.

For self-supervised learning, our proposed Masked Channel Modeling (MCM) technique is applied as described in section (IV). Specifically, 9 out of 64 real-part patches are randomly masked (approximately $15\%$ of patches), with the corresponding imaginary-part patches masked similarly to prevent data leakage. Masking is conducted based on three probability-based actions: with probability 0.1, patches are replaced with a random vector; with probability 0.8, they are replaced with a uniform MASK array (value $m = 1$); and with probability 0.1, they remain unchanged. This approach enables LWM to learn robust contextual relationships across both masked and unmasked patches, facilitating effective feature extraction from wireless channels. Importantly, because LWM does not know which patches are masked, it is compelled to understand all inter- and intra-patch relationships, rather than focusing solely on the masked patches. Although we only optimize the loss function based on the masked patches, this structure ensures that LWM captures a comprehensive representation across the entire input, reinforcing its ability to generalize effectively. The objective is to minimize the prediction loss (4) by selecting LWM embeddings $\mathbf{e}_i^{\text{LWM}}$ of the selected (masked) patches $\mathcal{M}$ and minimizing

$$\min_{\mathbf{W}^{\text{dec}}, \boldsymbol{\Theta}} \sum_{i \in \mathcal{M}} \left\| \mathbf{W}_i^{\text{dec}} \mathbf{e}_i^{\text{LWM}} - \mathbf{p}_i \right\|^2, \qquad (19)$$

where $\boldsymbol{\Theta}$ represents the parameters of the LWM model. This optimization encourages the model to learn mappings that accurately reconstruct masked patches from the embeddings, refining its ability to capture critical context from both masked and unmasked channel sections.

After masking, a CLS patch vector $\mathbf{p}_{\text{CLS}}$, initialized with random values, is prepended to the start of the channel patch sequence. Each $16 \times 1$ patch vector is then projected into a $D = 64$-dimensional space via a linear layer, yielding embeddings $\mathbf{e}_i^{\text{emb}} \in \mathbb{R}^{64}$ for each patch. These embeddings are further enriched with positional encodings to capture structural dependencies critical for learning within wireless channel data, as discussed in section (IV). The final input embedding matrix $\mathbf{E}^{\text{input}} \in \mathbb{R}^{129 \times 64}$ consolidates spatial and frequency information from both real and imaginary components, ready for processing through LWM's transformer-only encoder.

## C. Bidirectional Attention Mechanism

LWM employs a bidirectional attention mechanism, similar to models like BERT, where each patch in the input sequence attends to both previous and subsequent patches. This bidirectional structure enables the model to capture complex inter-patch dependencies across the channel matrix, making it more effective at understanding spatial and frequency relationships. Unlike autoregressive models like GPT, which process sequences in a strictly left-to-right (or right-to-left) manner, the bidirectional attention mechanism allows LWM to fully leverage the information from the entire channel when predicting masked patches. This is particularly important in wireless communications and sensing, where a subcarrier's channel might depend on both preceding and following subcarriers, and interactions between antennas must be understood in context.

Additionally, the transformer's self-attention mechanism, with its **multi-head attention**, plays a crucial role in this process by allowing the model to focus on different aspects of the channel simultaneously. For example, one attention head might focus on capturing local patterns, such as how signal strength varies across a few antennas, while another head might focus on larger-scale structures, such as the overall signal across different subcarriers. This enables the model to capture both fine-grained and broader patterns in wireless data, making it highly effective at extracting useful features.

## D. Pre-Training Loss Function

The pre-training of LWM uses the **MSE** loss function rather than the **cross-entropy** loss typically found in NLP models like BERT. In language models, cross-entropy loss is effective for predicting discrete tokens from a fixed vocabulary, a finite set of words, subwords, or characters that the model has been trained to recognize. Each word or token in the vocabulary is represented by a unique embedding, and during pre-training, the model learns to select the correct token by calculating probabilities across this vocabulary. For BERT, this is achieved through the Masked Language Modeling (MLM) task, where a certain percentage of tokens in the input sequence are masked, and the model is tasked with predicting the masked tokens based on the context provided by the remaining tokens. Since MLM is fundamentally a classification task, the goal is to identify the correct token from a pre-defined vocabulary, making cross-entropy loss suitable for assessing how closely the predicted probability distribution aligns with the true distribution for each masked token. The MLM loss function, designed to extract rich, contextual features from the surrounding text, is defined as

$$\mathcal{L}_{\mathsf{MLM}} = -\sum_{i \in \mathcal{M}} \log p(y_i | \mathbf{e}_i^{\mathsf{BERT}}), \quad (20)$$

where $\mathcal{M}$ denotes the set of masked positions, $y_i$ is the true token at position $i$, and $p(y_i | \mathbf{e}_i^{\mathsf{BERT}})$ is the predicted probability for $y_i$ based on the model's output embedding $\mathbf{e}_i^{\mathsf{BERT}}$ for the masked position. This enables the model to capture nuanced language representations, allowing a simple linear layer on top

to predict masked tokens by leveraging unmasked contextual information.

However, wireless channels differ fundamentally from discrete language data. Rather than discrete tokens, wireless channels consist of continuous, high-dimensional values representing complex spatial and temporal features. With no fixed vocabulary for continuous data, each patch can assume a broad range of values, making cross-entropy loss unsuitable. Instead, LWM employs MSE loss as in (4) for Masked Channel Modeling (MCM), which is essentially a regression task. MSE loss measures the error between predicted and actual continuous values of masked patches, guiding LWM to learn contextually rich features across the input sequence so that the linear layer can accurately predict masked patches based on the unmasked surrounding information. This task enables LWM to develop robust, meaningful feature representations tailored to the continuous nature of wireless channel data.

Table I: LWM Pre-training Setup Parameters

| Parameter | Value |
|---|---|
| Antennas at BS $N$ | 32 |
| Antennas at UEs | 1 |
| Subcarriers $M$ | 32 |
| Patch Size $L$ | 16 |
| Embedding Size $D$ | 64 |
| Channel Patches $P$ | 128 |
| Attention Heads $H$ | 12 |
| Encoder Layers $E$ | 12 |
| FFN Hidden Size $D_{\mathsf{FF}}$ | 256 |
| Head Dimension $D_{\mathsf{H}}$ | 5 |
| Learning Rate | $1 \times 10^{-4}$ |
| Batch Size | 64 |
| Optimizer | Adam |
| Adam $\beta_1$ | 0.9 |
| Adam $\beta_2$ | 0.999 |
| Adam $\epsilon$ | $1 \times 10^{-8}$ |
| Weight Decay | $1 \times 10^{-5}$ |
| Dropout Rate | 0.1 |
| Masking Percentage | 15% (80/10/10) |
| Model Parameters | 600K |
| Training Set Size | 820K |
| Validation Set Size | 200K |

## E. Pre-Training Setup Parameters

As shown in table I, the pre-training setup for LWM includes 12 attention heads, 12 encoder layers, an embedding size of 64, and an FFN hidden size of 256. Training begins with a learning rate of $1 \times 10^{-4}$, decreasing by 10% every 10 epochs to ensure smooth convergence. A batch size of 64 is utilized, along with the Adam optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.999$, eps $= 1 \times 10^{-8}$), and a weight decay of $1 \times 10^{-5}$ to reduce overfitting. The 12-head attention mechanism enables the model to capture multiple relationships within the data, while the depth of the encoder layers allows it to extract both

local and global patterns in the channels. This setup ensures effective learning, balancing convergence and generalization across various wireless communication scenarios.

By the end of pre-training, LWM is capable of producing rich, contextual embeddings from raw wireless channels. The integration of channel preprocessing, self-supervised learning, bidirectional attention, and multi-head attention mechanisms enables the model to generalize effectively across a range of scenarios, making it a powerful feature extractor for diverse downstream tasks in wireless communications and sensing systems.

## VII. INFERENCE

LWM excels in generating rich, context-aware embeddings from raw wireless channels in real-time, with no need for additional training for embedding generation. Pre-trained on large, diverse datasets using a self-supervised, Transformer-based approach, it allows users to immediately obtain high-quality, low- and high-dimensional embeddings suitable for a wide array of downstream tasks. The pre-trained model can be employed as-is, leveraging these embeddings directly, or choose to fine-tune the model's last layers to extract highly task-specific, fine-grained features. This flexibility makes LWM highly adaptable, enabling it to capture both local and global patterns and perform effectively across general and specialized scenarios, even in data-scarce environments.

The inference process begins with segmenting the raw wireless channel data into patches, embedding them, and adding positional encodings, similar to the pre-training phase but **without the need for masking or weight updates**. This structured approach allows LWM to extract and represent multi-scale patterns from both small and large contexts within the data. The resulting embeddings, $\mathbf{E}^{\text{LWM}} = [\mathbf{e}^{\text{LWM}}_{\text{CLS}}, \mathbf{e}^{\text{LWM}}_1, \mathbf{e}^{\text{LWM}}_2, \ldots, \mathbf{e}^{\text{LWM}}_P]^\top = \begin{bmatrix} \mathbf{C} & \mathbf{E}^\mathsf{T} \end{bmatrix}^\mathsf{T} \in \mathbb{R}^{(P+1) \times D}$, consist of the **CLS embedding** $\mathbf{C} \in \mathbb{R}^D$ and the **channel embeddings** $\mathbf{E} \in \mathbb{R}^{P \times D}$, as shown in Fig. 1.

The CLS embedding provides a holistic view of the channel, making it suitable for tasks that require general knowledge, such as LoS/NLoS classification. In contrast, the channel embeddings, which are four times larger than the input, capture complex spatial and frequency dependencies, making them ideal for more detailed applications requiring nuanced data insights. This inference setup offers multiple practical benefits:

- **Immediate Usability:** The pre-trained model generates embeddings that are ready for downstream use without further training, making LWM a plug-and-play solution for various applications. This efficiency proves invaluable in scenarios where resources are limited, and rapid, generalized feature extraction is essential for enhanced performance.
- **Multi-Scale Pattern Capture:** By transforming raw channels into both low- and high-dimensional embeddings, LWM captures local interactions within a patch as well as broader contextual relationships across patches. This multi-scale capability ensures that users can address

tasks requiring either fine-grained analysis or more generalized insights.
- **Task Flexibility:** The pre-trained embeddings can be applied directly or, if necessary, fine-tune the last layers of the model for enhanced task specificity. This flexibility is particularly useful in specialized applications that benefit from fine-grained features but do not necessarily require the entire model to be retrained.
- **Efficiency in Data-Scarce Environments:** LWM's ability to generalize effectively even with limited data proves advantageous in settings where large labeled datasets are unavailable. By leveraging the robust representations learned during pre-training, LWM maintains high performance with minimal data, making it well-suited to real-world conditions.

This structure enables LWM to excel in diverse wireless communication tasks by balancing detailed spatial-frequency capture and computational efficiency, providing a practical and adaptable inference process for a range of applications.

## VIII. DOWNSTREAM TASK EVALUATION

We demonstrate that LWM CLS and channel embeddings, generated in real-time, provides better performance for various downstream tasks, when compared to the original raw wireless channels.

### A. Sub-6 to mmWave Beam Prediction

This task aims to predict the strongest mmWave beam at the receiver from a predefined codebook at the base station, based on Sub-6 GHz channels. This is a specific case of *channel mapping*, where instead of directly predicting mmWave channels, the model learns the relationship between Sub-6 GHz channels and the best mmWave beam. Ground-truth optimal beams are computed for each user, generating a labeled dataset. The task complexity is varied by adjusting the codebook size, ranging from 16 to 256 beams. Performance is evaluated across different training data percentages to assess how well the model generalizes with less data. This makes the task highly practical for modern communication systems, as it reduces the overhead of full mmWave channel estimation. It also tests how well LWM embeddings capture the spatial and propagation characteristics of Sub-6 GHz channels and generalize to higher-frequency mmWave beams, making it an ideal benchmark for comparing the efficiency and generalizability of LWM embeddings against raw channels.

**Downstream Model**: To ensure fair comparison, we use a similar downstream model complexity for both embeddings and raw channels, selecting a model optimized for raw channels as the performance benchmark. A uniform Residual 1D-CNN architecture with 500K parameters is employed, designed to capture complex patterns through residual connections and weight-sharing while avoiding overfitting. The model consists of an initial convolution layer followed by three residual blocks, each containing several convolution layers that extract deeper features. This is followed by global average pooling and fully connected layers for final classification.

(a) Raw Channels Performance



(b) LWM Embeddings Performance



(c) Performance Difference
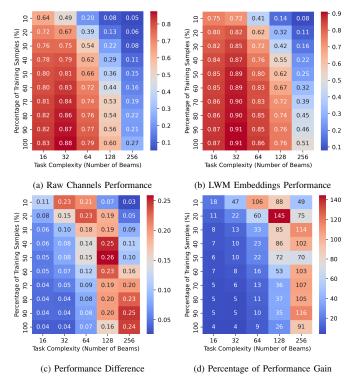


(d) Percentage of Performance Gain

Figure 3: This figure compares beam prediction F1-score performance between raw channels and their inferred LWM embeddings, based on a total of 10388 training raw channels, and highlights their relative effectiveness.

When parameters exceed 500K, the model overfits to raw channels, resulting in lower validation performance, which is why we consider this architecture the benchmark for raw channels.

**Downstream Task Dataset:** The test set comprises six new scenarios from the DeepMIMO dataset, which were not part of LWM's pre-training: *Denver*, *Fort Worth*, *Oklahoma*, *Indianapolis*, *Santa Clara*, and *San Diego*, comprising a total of 14840 samples. For each user, the corresponding 28GHz channels are generated, and the optimal beams are computed to create a labeled dataset, where 3.5GHz channels serve as inputs and 28GHz beams as labels. The dataset is split into 70% for training, 20% for validation, and 10% for testing. For LWM embeddings, the raw channels are first to generate LWM embeddings, the raw channels are processed through the pre-trained LWM model in real-time, with the channel embedding component $\mathbf{E}$ serving as the input for this downstream model, as shown in Fig. 1.

**Downstream Tasks Evaluation:** Fig. 3 presents a comparison of the performance between raw channels and LWM channel embeddings in training a model for beam prediction, evaluated across different codebook sizes and varying amounts of training data. As shown in Fig. 3a, the downstream models trained with raw channels require significantly more data to reach high performance levels, while LWM channel embeddings usually achieve performance saturation with just 50%-70% of the available data and consistently outperform raw channels. Notably, LWM embeddings reach the benchmark

performance of raw channels with only 40%-50% of the data, regardless of task complexity, highlighting the data efficiency of LWM embeddings.

We use the F1-score for classification tasks as it accounts for imbalanced labels and provides a clearer evaluation of model performance than accuracy. The F1-score difference heatmap in Fig. 3c highlights where embeddings surpass raw channels, aiding in model optimization for complex or low-data tasks. Meanwhile, the F1-score gain percentage heatmap in Fig. 3d emphasizes the efficiency of embeddings, showing how they scale with fewer resources. Fig. 3c demonstrates that embeddings are most effective when data is limited relative to task complexity. The performance difference follows a concave trend, where an initial increase in data improves performance, but beyond a certain point, the difference diminishes and gradually decreases, though it never turns negative.

As the task complexity increases, such as with a higher number of labels, larger datasets are required to fully showcase the advantages of embeddings over raw channels, allowing embeddings to reveal their full potential in capturing intricate patterns and dependencies.

### B. LoS/NLoS Classification

This task serves as a great benchmark to evaluate CLS embeddings, which provide a highly compressed but more informative representation of raw channels. These embeddings effectively capture the critical features of a channel, making them ideal for various tasks that require a holistic understanding of the channel's behavior. For instance, in applications like UE CSI feedback, the CLS embeddings can reduce the overhead of sending full channel data to the base station, particularly in multi-vendor environments. Instead of relying on infeasible joint autoencoder training between user equipment (UE) and base station (BS), the LWM model functions like an encoder, delivering a compact yet illustrative version of the channel. This makes it highly suitable for tasks that need a comprehensive understanding of channel characteristics, such as beamforming or interference management, while significantly reducing transmission complexity.

Since this is a relatively simple task, we utilize a fully-connected network (FCN) with 30K parameters to classify LoS and NLoS channels. To ensure fairness, we use a deeper model for the CLS embeddings, keeping the total number of parameters the same. The same six scenarios are used for this task, with the labels changed to binary values indicating LoS or NLoS channels. During evaluation, the smaller CLS embeddings significantly outperform raw channels, achieving an F1-score of 0.87 with just 13 training samples, compared to 0.55 for raw channels, highlighting the superior generalization of the embeddings. Although the performance gap reduces with more training data, raw channels never exceed the performance of the embeddings.

### C. Robust Beamforming

Beamforming design is relatively straightforward when the accurate channel information is available. However, in practi-

cal scenarios, errors in the estimated channel can significantly degrade beamforming performance, leading to suboptimal signal-to-noise ratios (SNR) and reduced reliability. LWM can be applied to mitigate these errors by implicitly denoising the estimated channel through robust feature extraction, enhancing performance even when parts of the channel information are missing or noisy.

The robust beamforming task is ideal for evaluating LWM embeddings due to the need for capturing intricate spatial and frequency relationships in wireless channels to predict optimal beam configurations. LWM's detailed feature extraction capabilities are especially well-suited for this task, showcasing its effectiveness in handling regression tasks beyond classification. The MCM pre-training approach enhances the model's robustness to noisy or imperfect CSI, making LWM embeddings practical for scenarios where complete or accurate channel information is unavailable.

We formulate an SNR maximization problem for users under a power constraint at the BS and solve it using the maximum ratio transmission (MRT) method. We then generate a dataset from this setup and train a lightweight downstream FCN with 30K parameters, mapping channels to beamforming vectors at the BS. Notably, when test channels contain masked patches (simulating imperfect CSI), LWM embeddings maintain robust performance, achieving an MSE of 0.01, compared to 0.51 with raw channels. This substantial improvement highlights the model's robustness and adaptability, demonstrating LWM's utility in both classification and regression tasks.

## IX. Discussion and Remarks

In this section, we provide additional insights related to the LWM framework, covering points on channel segmentation, attention mechanisms, Masked Channel Modeling (MCM) pre-training, and a comparison with autoencoder embeddings. These discussions offer complementary perspectives to enhance understanding of the model's design and functionality.

### A. Discussion on Channel Segmentation

Dividing channels into patches, rather than using each element individually, provides several advantages in terms of computational efficiency, contextual learning, and the ability to capture both local and global relationships within data sequences. Below, we outline these advantages.

**Computational Efficiency and Contextual Learning:** When each data element is processed individually, the model might capture only fine-grained, isolated information without effectively understanding broader spatial or temporal structures. This can lead to inefficiencies in terms of computational load and a limited capacity for extracting meaningful context, especially when working with data that has rich, interdependent structures.

**Hierarchical Structure for Enhanced Representation:** By grouping data elements into patches, we introduce a hierarchical structure that allows the model to treat each patch as a higher-level feature, capturing a more integrated view of the data. Patches inherently represent small local contexts, which

may correspond to subsets of adjacent measurements, neighboring time points, or correlated data points. This grouping enables the model to interpret these elements together, learning intra-patch dependencies that reflect localized patterns and reduce noise from treating elements in isolation. Consequently, the patches serve as a natural bridge between the fine-grained details and the broader context, with each patch encompassing more significant information than individual points would.

**Capturing Short-Range and Long-Range Dependencies:** When patches are treated as individual tokens in the Transformer model, self-attention mechanisms can efficiently learn inter-patch relationships, establishing contextual dependencies across patches. This is crucial for capturing both short-range and long-range dependencies within the data, as the self-attention mechanism can evaluate how each patch relates to every other patch. Such inter-patch dependencies are challenging to capture when processing individual elements alone, as the model would lack the structured grouping needed to define higher-order relationships effectively.

**Similarity to CNNs in Capturing Spatial Relationships:** Patch-based self-attention and CNNs share a similarity in capturing spatial relationships, though their methods differ. In CNNs, convolutional filters slide across the input in a fixed grid, capturing local patterns with each filter pass. This process enables CNNs to identify localized features, like edges or textures, through layers that gradually increase the receptive field, allowing the model to build on local patterns to understand larger, more complex structures. Patch-based self-attention, though, can achieve a similar effect by treating patches as analogous to these localized features, effectively capturing the relationships within each patch while simultaneously considering inter-patch relationships. This approach offers the benefits of a CNN's hierarchical feature extraction, with each patch resembling a "receptive field" in a convolution layer, yet goes further by allowing direct, learnable interactions between distant patches in one layer. Thus, patch-based self-attention maintains the CNN advantage of learning spatial hierarchies but improves upon it by enabling a global view without stacking additional layers, making it more versatile for complex data patterns.

**Alignment with Multi-Head Attention Mechanism:** Another benefit of the patch-based approach is that it aligns more naturally with the multi-head attention mechanism in Transformers. Each head can focus on different parts of the data or specific features within the patches, creating diverse representations that can capture various aspects of the underlying structure in the data. This parallel processing across patches not only accelerates computations but also provides a more nuanced view of the data, enabling the model to leverage both focused and distributed patterns that contribute to robust feature extraction.

**Enhanced Masked Channel Modeling (MCM):** The patch-based approach also improves the effectiveness of our proposed MCM self-supervised learning approach by introducing a more challenging prediction task: when an entire patch is masked, the model must predict it without access to any of

its internal neighboring elements. This requires the model to leverage information from other patches, encouraging a deeper understanding of global relationships within the data. As a result, the patch-based approach leads to richer, more context-aware embeddings that are highly effective for downstream tasks.

**Flexibility in Model Design:** Finally, the patch-based approach allows for flexible design choices in the model architecture. By controlling the patch size, we can tune the model to balance fine detail with computational efficiency, adjusting the representation granularity as needed. Large patches reduce the sequence length, making it more computationally feasible to work with extensive datasets, while small patches retain more fine-grained detail at the cost of higher computation. This flexibility, combined with the capacity to capture multi-scale dependencies, makes patch-based data preparation a powerful method for extracting rich and contextually meaningful embeddings from complex data.

### B. Discussion on Attention

LWM's Transformer-based architecture consists of $E$ stacked encoder layers, each layer comprising a self-attention sub-layer and an FFN sub-layer, with layer normalization and residual connections applied after each. This layered configuration provides hierarchical abstraction, progressively refining the representation of the input patches. The Transformer encoder, through parallel attention mechanisms and deep stackable layers, becomes highly efficient for extracting features from wireless channels, ensuring it captures a diverse set of dependencies. Thus, the model's architecture is robust, allowing it to process data with spatial and temporal patterns crucial in wireless communication and sensing scenarios.

The self-attention mechanism in Transformer models allows each element in a data sequence to attend to all other elements, a powerful capability that enables the model to capture essential patterns, dependencies, and multi-faceted perspectives within the data. By computing contextually weighted representations, self-attention mechanisms extract and emphasize relationships across data sequences in a mathematically optimal way. We explore several practical examples below.

**Identifying and Prioritizing Significant Data Patterns:** For data sequences where only a subset of points significantly affects the desired output, attention scores can isolate and prioritize the most relevant patterns. Given a sequence of data points $\{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_N\}$, attention assigns a relevance score to each pair of points. Specifically, each data point is represented by a query $\mathbf{q}_i$, key $\mathbf{k}_j$, and value $\mathbf{v}_j$ vector, with relevance computed as follows

$$\alpha_{ij} = \frac{\exp\left(\frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{D}}\right)}{\sum_{j=1}^{N} \exp\left(\frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{D}}\right)}, \tag{21}$$

where $\alpha_{ij}$ is the attention weight, or relevance score, which emphasizes pairs $(i, j)$ with a high query-key similarity. For each query, the output is a weighted sum of the values

$$\mathbf{z}_i = \sum_{j=1}^{N} \alpha_{ij} \mathbf{v}_j. \tag{22}$$

The score $\mathbf{z}_i$ essentially filters out less relevant points, efficiently highlighting critical data patterns by amplifying highly weighted points while diminishing the impact of lower-weighted ones.

**Capturing Long-Range Dependencies in Sequential Data:** When modeling sequential data, long-range dependencies, where distant elements influence each other, are crucial. Self-attention is particularly effective here, as each point $\mathbf{x}_i$ can attend to every other point in the sequence, capturing dependencies regardless of distance. With queries, keys, and values represented as $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times D}$, respectively, the attention score matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ is calculated as follows

$$\mathbf{A} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right), \tag{23}$$

where each entry $A_{ij}$ captures the influence between elements $\mathbf{x}_i$ and $\mathbf{x}_j$. Thus, the final output matrix is computed as

$$\mathbf{Z} = \mathbf{A}\mathbf{V}, \tag{24}$$

where distant dependencies are inherently represented in $\mathbf{Z}$ without explicitly increasing the computational complexity. This facilitates capturing long-range relationships in the data, which is vital for any sequence-based processing.

**Aggregating Multi-Faceted Perspectives of Data with Multi-Head Attention:** In cases where the same data element may have multiple interpretations, multi-head attention allows the model to process these interpretations in parallel. Each head $h$ is an independent attention mechanism with separate learned parameters for query, key, and value projections, denoted by $\mathbf{W}_h^Q, \mathbf{W}_h^K, \mathbf{W}_h^V \in \mathbb{R}^{D \times D_\mathsf{H}}$. The output for each head is calculated as follows

$$\text{head}_h = \text{softmax}\left(\frac{(\mathbf{Q}\mathbf{W}_h^Q)(\mathbf{K}\mathbf{W}_h^K)^\top}{\sqrt{D_\mathsf{H}}}\right)(\mathbf{V}\mathbf{W}_h^V). \tag{25}$$

These head outputs are concatenated and transformed through a linear weight matrix $\mathbf{W}^O \in \mathbb{R}^{D_H H \times D}$, giving

$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\text{head}_1, \ldots, \text{head}_H)\mathbf{W}^O. \tag{26}$$

This approach provides multiple "views" of each input, allowing the model to capture different relationships for the same data point. These multi-faceted representations are particularly useful in complex data sets where different perspectives on a single data point could yield unique insights.

**Dynamic Processing Based on Contextual Relationships:** Self-attention dynamically adjusts each data element's representation based on its context within the sequence. For each element $\mathbf{x}_i$, the similarity between its query $\mathbf{q}_i$ and all keys

**K** determines the relevance of every other element in the sequence, yielding a contextually weighted output

$$\mathbf{z}_i = \sum_{j=1}^{N} \alpha_{ij} \mathbf{v}_j. \tag{27}$$

This adaptive approach, guided by the context provided by neighboring elements, allows the model to encode each point's relationships within the broader sequence, resulting in outputs $\mathbf{z}_i$ that dynamically adapt based on the sequence structure. This adaptability is key for understanding complex, interdependent data.

**Interpretability of Data Relationships Through Attention Scores:** A notable byproduct of self-attention is its interpretability. Attention scores $\alpha_{ij}$ provide insight into which data elements contribute most to each output, as each score represents the relative importance of $\mathbf{x}_j$ to $\mathbf{x}_i$. This interpretability can be quantified through the attention matrix

$$\mathbf{A} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{D}}\right). \tag{28}$$

By analyzing **A**, one can understand how each data element influences the output, which is essential for tasks requiring explainability and insight into model decision-making. This feature makes self-attention-based architectures valuable for applications where understanding input-output relationships is as important as achieving high accuracy.

### C. Discussion on Masked Channel Modeling

Masked Channel Modeling (MCM) is a critical pre-training strategy for LWM that enhances its robustness to noise. By masking a subset of patches (in this case, $15\%$) and training the model to reconstruct them accurately, MCM encourages LWM to develop a deep understanding of the underlying structure and dependencies within the channel data. This process enables the model to infer missing or corrupted information by leveraging context from unmasked patches. Consequently, the pre-trained LWM becomes resilient to noise and distortions in channel data, as it learns to rely on contextual information rather than exact values alone.

During inference, this robustness manifests as improved performance when the model encounters noisy channel inputs. Because the model has been trained to reconstruct masked portions, it effectively learns a form of **noise-agnostic feature representation**, allowing it to handle noisy or incomplete data inputs without significant performance degradation. This makes MCM an effective approach for achieving noise robustness in wireless communication and sensing applications.

Beyond robustness, the MCM-pretrained LWM can also serve directly as a channel denoiser during inference. In this application, noisy channels are input to the model, and the LWM's learned embeddings work to produce a refined representation of the original signal. The model's linear layer, trained to map embeddings back to the original patch dimensions, effectively reconstructs patches in a way that minimizes the impact of noise. The MSE minimization objective during pre-training allows the model to prioritize accurate reconstruction of essential signal components over noise artifacts. Therefore, in practice, LWM can filter out noise by reconstructing the signal patches based on learned dependencies and patterns. This makes LWM a valuable tool for channel denoising, improving overall signal quality and communication reliability in wireless environments where noise is a significant concern.

Through MCM, LWM gains both robustness to noise and denoising capability, making it versatile for various inference applications in imperfect channels.

### D. LWM vs. Autoencoders Embeddings

LWM and autoencoders, while both designed for feature extraction, represent fundamentally different approaches in terms of architecture, training objectives, and applicability to downstream tasks. Below, we explore these differences from several perspectives.

**Architectural Differences:** Autoencoders consist of two main components: an encoder $f_{\text{enc}} : \mathbb{R}^n \to \mathbb{R}^d$ that compresses input data $\mathbf{x} \in \mathbb{R}^n$ into a lower-dimensional latent space representation $\mathbf{z} \in \mathbb{R}^d$ (with $d < n$), and a decoder $f_{\text{dec}} : \mathbb{R}^d \to \mathbb{R}^n$ that attempts to reconstruct the original input from $\mathbf{z}$. The training objective minimizes a reconstruction loss

$$\mathcal{L}_{\text{AE}} = \frac{1}{N} \sum_{i=1}^{N} \|\mathbf{x}_i - f_{\text{dec}}(f_{\text{enc}}(\mathbf{x}_i))\|^2, \tag{29}$$

where $N$ is the number of training samples. This setup encourages the encoder to learn features that capture the input's primary structure but focuses on elements necessary for reconstruction rather than on generalized representations for other tasks.

LWM, on the other hand, uses a Transformer-based architecture designed to learn contextual embeddings through self-attention rather than reconstruction. Given a sequence of input patches $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_P\}$, where each patch represents a portion of the data, the self-attention mechanism generates output embeddings **E** by attending to each patch in the sequence based on its relevance to others. The attention weights $\alpha_{ij}$ are calculated as follows

$$\alpha_{ij} = \frac{\exp\left(\mathbf{q}_i \cdot \mathbf{k}_j / \sqrt{D}\right)}{\sum_{j=1}^{P} \exp\left(\mathbf{q}_i \cdot \mathbf{k}_j / \sqrt{D}\right)}, \tag{30}$$

where $\mathbf{q}_i$ and $\mathbf{k}_j$ are the query and key vectors for patches $i$ and $j$, respectively. This approach allows LWM to capture relationships between patches, yielding embeddings enriched with local and global contextual information.

**Training Objective and Generalization:** Autoencoders are generally trained to minimize reconstruction error, which encourages the model to compress information in a way that maximally preserves input data for decoding. While this is effective for dimensionality reduction, it limits the model's ability to capture task-specific features unless specifically fine-tuned for each task. Additionally, reconstruction-oriented training often struggles with generalization since the focus

is on faithfully representing the input rather than capturing diverse, adaptable representations.

In contrast, LWM employs MCM, a self-supervised learning objective that masks certain patches and tasks the model with predicting these masked sections. For a set of masked patches $\mathcal{M}$, the objective is to minimize

$$\mathcal{L}_{\mathsf{MCM}} = \frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \left\| \mathbf{W}_i^{\mathsf{dec}} \mathbf{e}_i^{\mathsf{LWM}} - \mathbf{p}_i \right\|^2, \tag{31}$$

where $\mathbf{e}_i^{\mathsf{LWM}}$ is the embedding for masked patch $i$, $\mathbf{W}_i^{\mathsf{dec}}$ is a linear mapping from the embedding space to the original patch space, and $\mathbf{p}_i$ is the original patch value. MCM forces LWM to learn not only local features within each patch but also dependencies across patches, yielding a richer, context-aware embedding that generalizes well across tasks.

**Handling Context and Dependencies:** In autoencoders, each input element is generally processed independently by the encoder, and contextual relationships between elements must be encoded indirectly through the encoder's weights. This is challenging for data with complex dependencies, as these models lack a mechanism to explicitly capture interactions between spatial or sequential elements.

LWM, through self-attention, can explicitly model these dependencies. For example, if each patch $\mathbf{x}_i$ represents a localized region of a larger data structure, the attention mechanism allows LWM to dynamically weigh the importance of other patches when forming each patch's representation. This property enables LWM to capture both local dependencies (similar to the receptive field in CNNs) and long-range dependencies by attending to distant patches, which autoencoders struggle to achieve efficiently.

**Efficiency in Real-Time Inference and Adaptability:** Autoencoders require both encoding and decoding steps, which can increase inference time, especially if fine-grained reconstruction is necessary. Moreover, adapting autoencoders to new tasks often requires re-training or fine-tuning, limiting flexibility.

In LWM, the embeddings generated during inference are immediately task-ready. There is no need for decoding, as the model's output is a high-quality embedding that can be directly used for downstream tasks. The self-attention mechanism, despite being computationally intensive, is highly parallelizable, allowing for efficient real-time inference with large datasets. Additionally, since LWM is pre-trained with a generalized objective, its embeddings are highly adaptable, capable of performing well across a wide range of tasks without additional training.

**Flexibility in Representational Capacity:** Autoencoders typically condense information into a fixed-size latent vector $\mathbf{z}$ that represents the entire input. This approach limits flexibility, especially for complex inputs where the representational needs may vary across sections. In LWM, each patch is embedded independently, allowing for both high-dimensional representations (capturing detailed features) and lower-dimensional summaries (e.g., CLS patches) for more generalized tasks.

This adaptive representational capacity makes LWM a better fit for scenarios that require a dynamic range of feature extraction, from detailed spatial features to broader contextual summaries.

To conclude, while autoencoders offer straightforward dimensionality reduction through reconstruction, LWM provides a flexible, context-aware representation by leveraging self-attention and self-supervised learning. These design choices allow LWM to capture complex dependencies and generalize effectively to diverse downstream tasks, making it a powerful alternative to autoencoders in contexts that demand rich, adaptable feature extraction.

## X. CONCLUSION

In this work, we introduced LWM, a pre-trained foundation model specifically designed for wireless communication and sensing channels. LWM draws inspiration from LLMs and Vision Transformers. Using a masked channel modeling strategy within a self-supervised framework, LWM is pre-trained to predict masked portions of the input via self-attention features and simple linear layers. Once pre-trained, LWM can generate real-time embeddings for raw channels, extracting rich, complex patterns from the input and consistently outperforming raw channels in downstream tasks. The LWM model is publicly available now.

## REFERENCES

[1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," 2023.

[2] T. S. Rappaport, Y. Xing, O. Kanhere, S. Ju, A. Madanayake, S. Mandal, A. Alkhateeb, and G. C. Trichopoulos, "Wireless communications and applications above 100 GHz: Opportunities and challenges for 6G and beyond," *IEEE access*, vol. 7, pp. 78729–78757, 2019.

[3] W. Jiang, B. Han, M. A. Habibi, and H. D. Schotten, "The road towards 6G: A comprehensive survey," *IEEE Open Journal of the Communications Society*, vol. 2, pp. 334–366, 2021.

[4] Z. Zhang, Y. Xiao, Z. Ma, M. Xiao, Z. Ding, X. Lei, G. K. Karagiannidis, and P. Fan, "6G wireless networks: Vision, requirements, architecture, and key technologies," *IEEE Vehicular Technology Magazine*, vol. 14, no. 3, pp. 28–41, 2019.

[5] 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.

[6] W. Yu, F. Sohrabi, and T. Jiang, "Role of deep learning in wireless communications," *IEEE BITS the Information Theory Magazine*, vol. 2, no. 2, pp. 56–72, 2022.

[7] M. Alrabeiah and A. Alkhateeb, "Deep learning for mmWave beam and blockage prediction using sub-6 GHz channels," *IEEE Transactions on Communications*, vol. 68, no. 9, pp. 5504–5518, 2020.

[8] M. Alrabeiah and A. Alkhateeb, "Deep learning for TDD and FDD massive MIMO: Mapping channels in space and frequency," in *2019 53rd asilomar conference on signals, systems, and computers*, pp. 1465–1470, IEEE, 2019.

[9] A. Jagannath, J. Jagannath, and T. Melodia, "Redefining wireless communication for 6g: Signal processing meets deep learning with deep unfolding," *IEEE Transactions on Artificial Intelligence*, vol. 2, no. 6, pp. 528–536, 2021.

[10] A. Diba, M. Fayyaz, V. Sharma, A. H. Karami, M. M. Arzani, R. Yousefzadeh, and L. Van Gool, "Temporal 3D convnets: New architecture and transfer learning for video classification," *arXiv preprint arXiv:1711.08200*, 2017.

[11] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2625–2634, 2015.

[12] D. Bahdanau, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.

[13] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," 2019.

[14] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. d. l. Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, *et al.*, "Mistral 7B," *arXiv preprint arXiv:2310.06825*, 2023.

[15] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, *et al.*, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.

[16] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, *et al.*, "GPT-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.

[17] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in neural information processing systems*, vol. 33, pp. 12449–12460, 2020.

[18] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *International conference on machine learning*, pp. 28492–28518, PMLR, 2023.

[19] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM transactions on audio, speech, and language processing*, vol. 29, pp. 3451–3460, 2021.

[20] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," 2021.

[21] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan, and M. Shah, "Transformers in vision: A survey," *ACM Computing Surveys*, vol. 54, p. 1–41, Jan. 2022.

[22] K. Han, Y. Wang, H. Chen, X. Chen, J. Guo, Z. Liu, Y. Tang, A. Xiao, C. Xu, Y. Xu, *et al.*, "A survey on vision transformer," *IEEE transactions on pattern analysis and machine intelligence*, vol. 45, no. 1, pp. 87–110, 2022.

[23] C.-F. Chen, Q. Fan, and R. Panda, "Crossvit: Cross-attention multi-scale vision transformer for image classification," 2021.

[24] L. Wu, W. Zhang, T. Jiang, W. Yang, X. Jin, and W. Zeng, "[cls] token is all you need for zero-shot semantic segmentation," 2023.

[25] A. Alkhateeb, "DeepMIMO: A generic deep learning dataset for millimeter wave and massive MIMO applications," in *Proc. of Information Theory and Applications Workshop (ITA)*, (San Diego, CA), pp. 1–8, Feb 2019.