

Here is a structured set of bullet points summarizing the main points of the content as if creating slides for a presentation:

#### **Slide 1: Introduction**

- Foundation models (FMs) are large models pretrained on massive data then adapted for downstream tasks
- Sequence models are the backbone of FMs, operating on arbitrary sequences from various domains

#### **Slide 2: Challenges with Current Models**

- Modern FMs are predominantly based on a single type of sequence model: Transformer and its core attention layer
- The efficacy of self-attention is attributed to its ability to route information densely within a context window
- However, this property brings fundamental drawbacks, such as inability to model outside a finite window and quadratic scaling with respect to the window length

#### **Slide 3: Recent Advances in Sequence Modeling**

- Structured state space sequence models (SSMs) have emerged as a promising class of architectures for sequence modeling
- SSMs can be interpreted as a combination of recurrent neural networks (RNNs) and convolutional neural networks (CNNs)

#### **Slide 4: Benefits of SSMs**

- Computationally efficient, with linear or near-linear scaling in sequence length
- Can be computed very efficiently as either a recurrence or convolution

#### **Slide 5: Conclusion**

- Mamba achieves state-of-the-art performance across several modalities, including language, audio, and genomics
- Mamba-3B model outperforms Transformers of the same size and matches its size twice in pretraining and downstream evaluation.