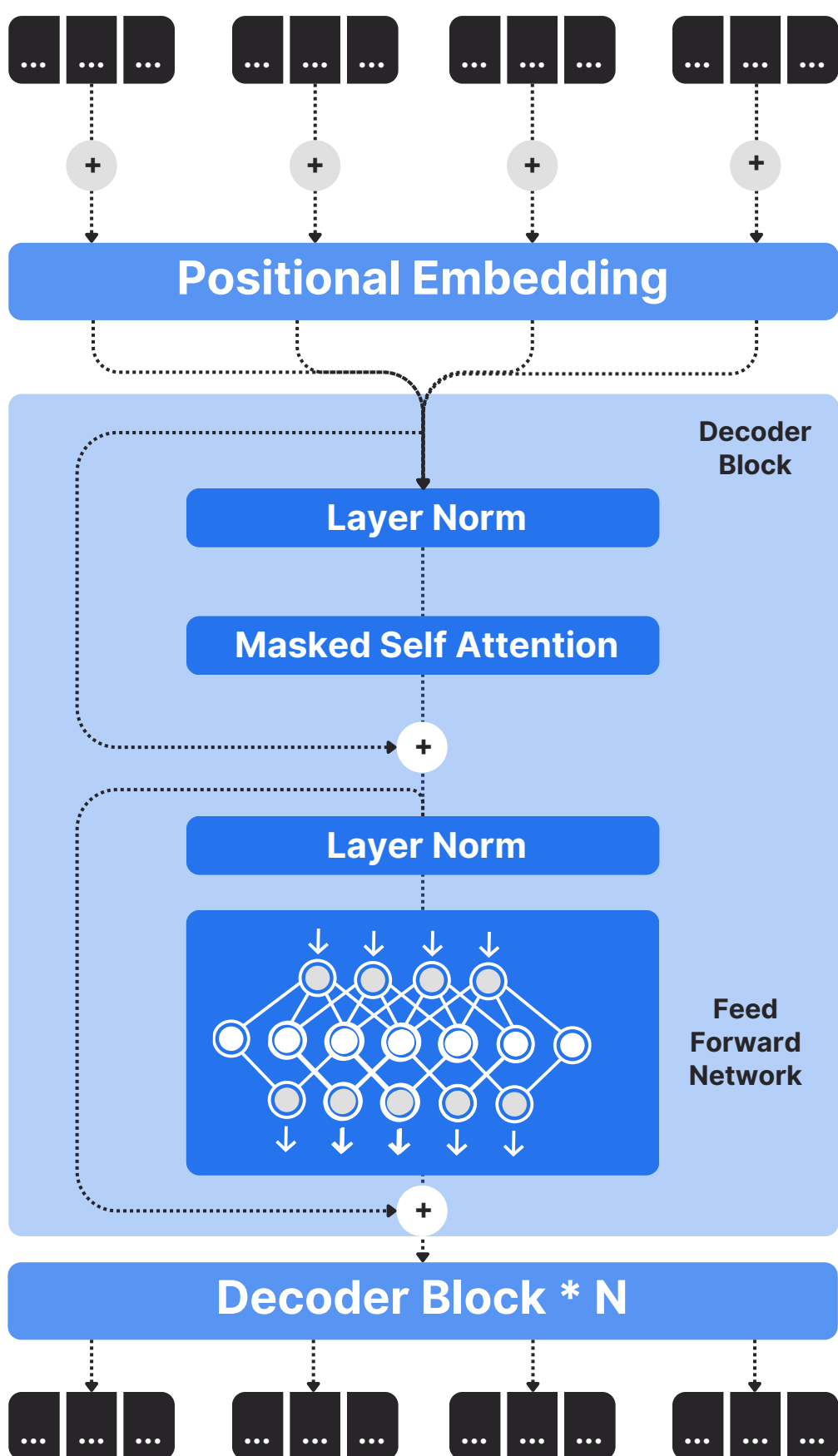
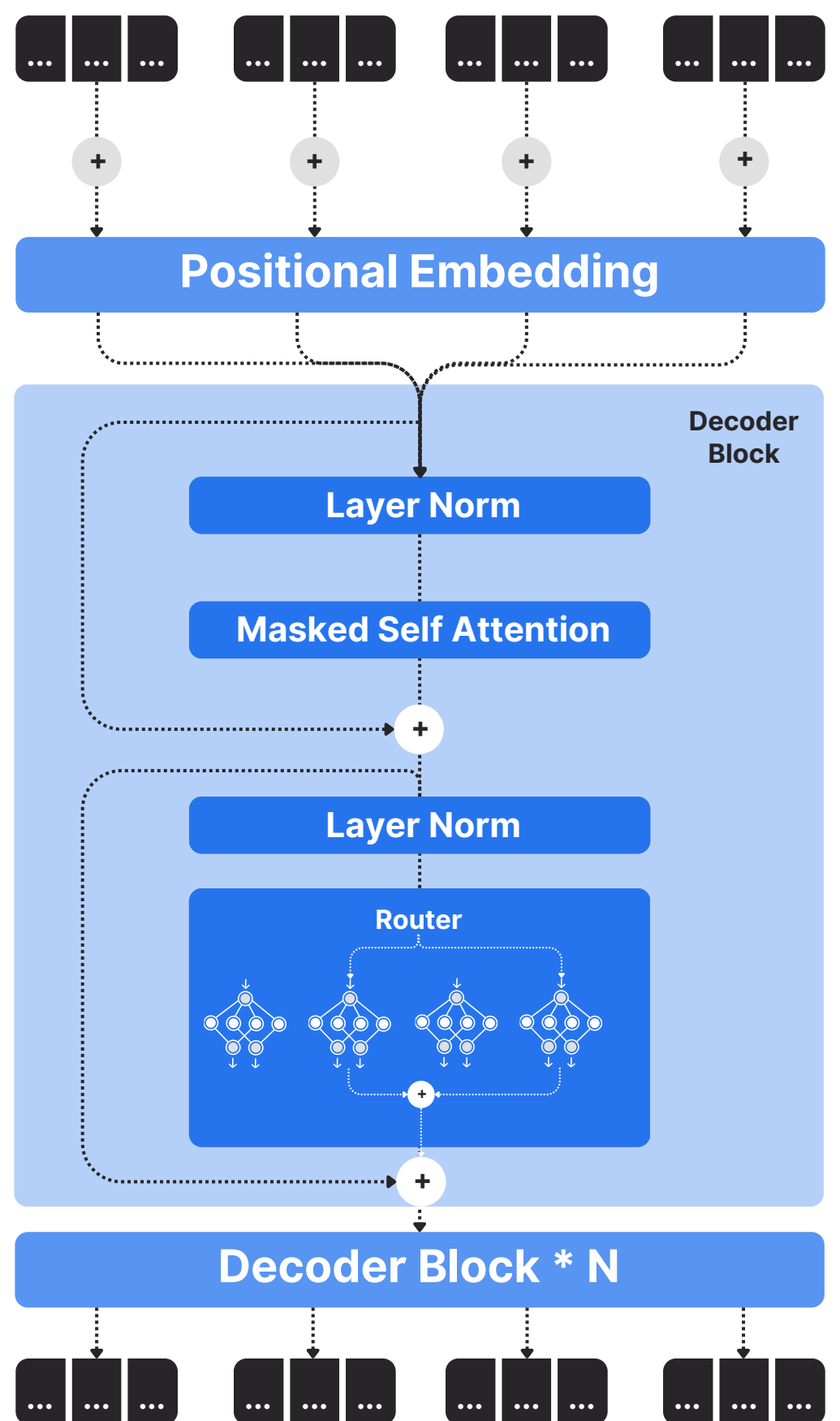


Transformers vs MoE

Transformers

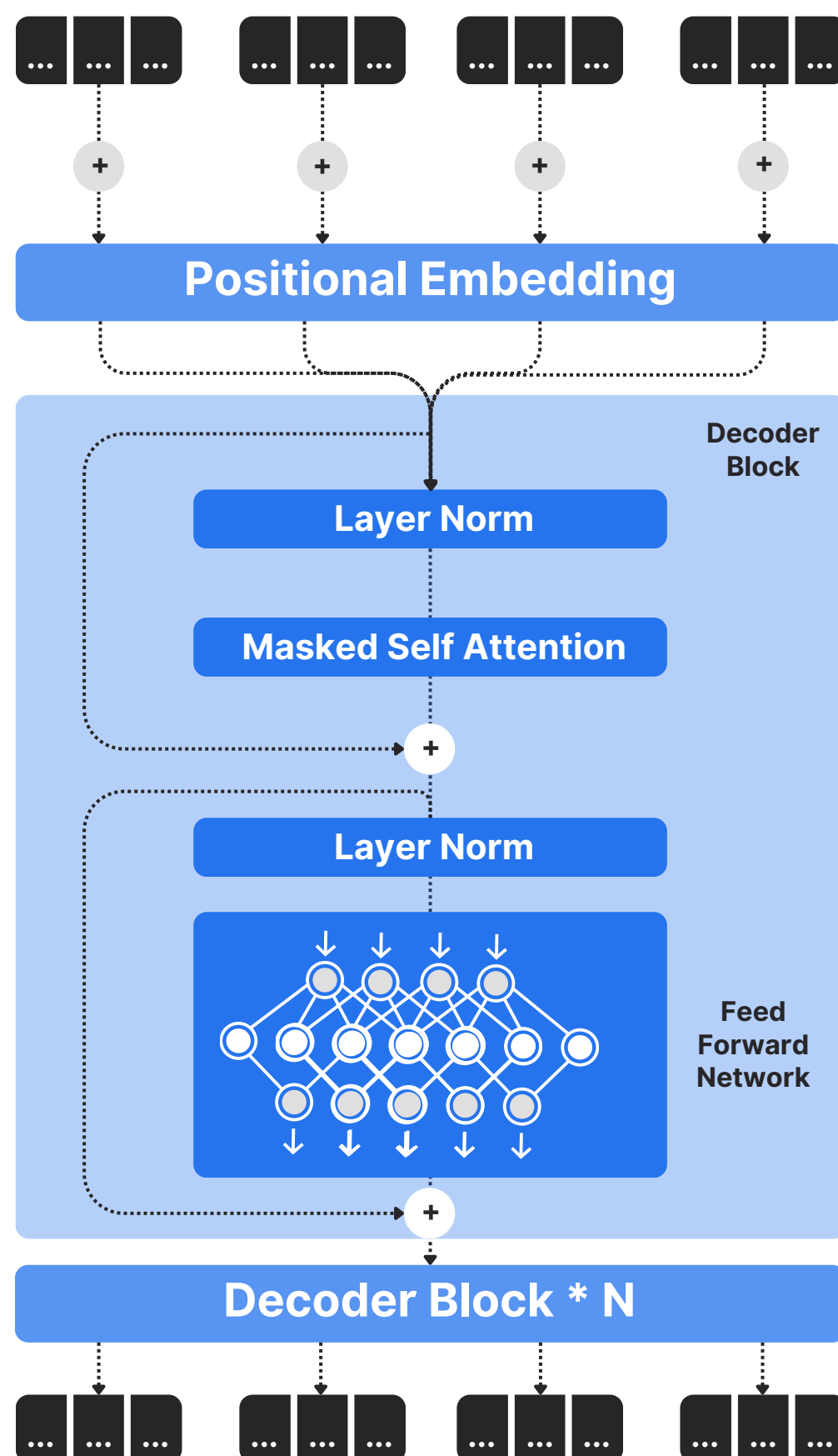


Mixture of Experts



What is a Transformer?

- The Transformer architecture, introduced in the paper "Attention is All You Need" (Vaswani et al., 2017), is a deep learning model that relies on the self-attention mechanism and positional encoding to process sequential data efficiently.



Key Components

- **Self-Attention Mechanism:** Enables the model to weigh the importance of different words in a sequence relative to one another.
- **Multi-Head Attention:** Improves the ability to capture different contextual dependencies.
- **Feedforward Layers:** Applies transformations to each token independently after attention computation.
- **Positional Encoding:** Injects order information into the model since Transformers do not have inherent sequential bias.
- **Layer Normalization and Residual Connections:** Ensure stable training and deeper architectures.

Advantages of Transformers

- **Parallelization:** Unlike recurrent models (RNNs, LSTMs), Transformers process sequences in parallel, leading to efficient training.
- **Scalability:** Can handle large datasets effectively (e.g., GPT, BERT, T5).
- **State-of-the-art Performance:** Achieves superior results in NLP, vision, and multimodal tasks.

What is MoE?

- Mixture of Experts (MoE) is a neural network architecture that dynamically selects a subset of specialized sub-models ("experts") to process each input. This approach improves efficiency by activating only relevant experts rather than using the full model for every input.

Key Components

- **Experts:** Individual neural network sub-models, each trained to specialize in a particular subset of data.
- **Gating Network:** A trainable component that determines which experts to activate for a given input.
- **Sparse Activation:** Unlike Transformers, which fully activate all layers, MoE selectively activates only a few experts per inference step, leading to computational efficiency.

Transformer vs. MoE: Key Differences

Feature	Transformer	Mixture of Experts (MoE)
Computation Type	Fully dense computation	Sparse computation (only a few experts activate per step)
Efficiency	Computationally expensive as all parameters are always active	More efficient by selectively using experts
Scalability	Scales with compute but requires full model activation	Can scale massively while keeping computation constant per token
Specialization	General-purpose	Experts can specialize in sub-tasks or domains
Latency	High due to dense computation	Lower since fewer parameters are used per forward pass
Model Complexity	High, with fully connected layers and self-attention	Higher due to gating networks and multiple experts

Recent architectures like Switch Transformers (Fedus et al., 2021) integrate MoE within Transformer layers, allowing large-scale training with significantly reduced computation. Key innovations include:

- Sparse Gated Layers within Transformers, replacing dense feedforward layers.
- Load balancing mechanisms to ensure fair expert usage.
- Improved training stability using routing strategies.
- This hybrid approach combines the best of both worlds: the expressive power of Transformers and the efficiency of MoE.

Use Cases

Application	Transformer	MoE
NLP (Chatbots, Translation)	✓	✓
Large-Scale Language Models	✓	✓ (GPT-4, GLaM)
Computer Vision	✓ (ViTs)	Limited
Multimodal Learning	✓	✓
Edge AI (Low Compute)	× (High cost)	✓ (Sparsity helps)

Both Transformers and Mixture of Experts have their strengths and trade-offs. While Transformers are powerful in handling sequential data and offer parallel computation, MoE provides a scalable approach by distributing computation across specialized experts.

The combination of these two architectures, as seen in Switch Transformers and GLaM, is paving the way for even more efficient and powerful AI models.