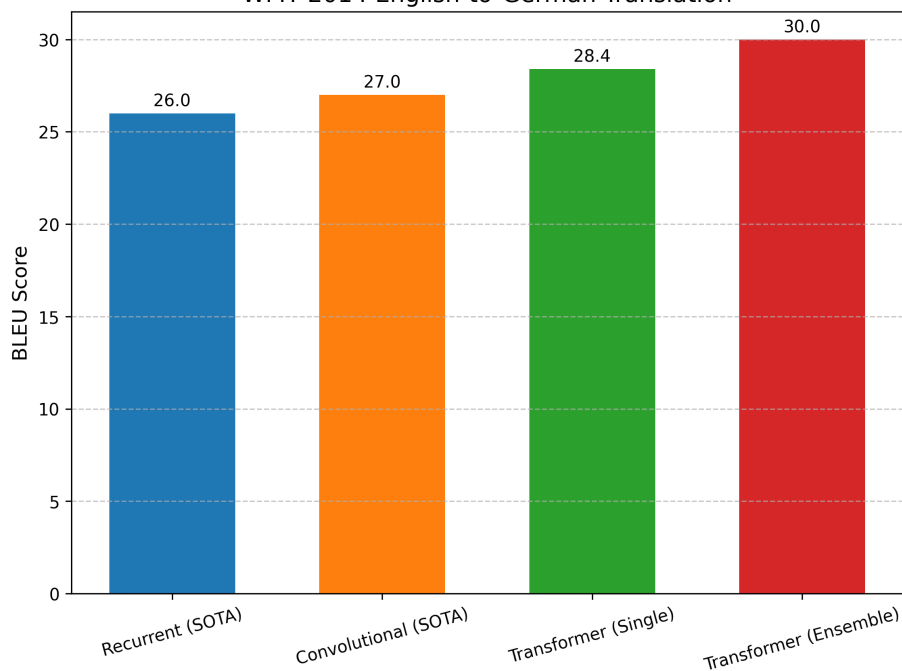
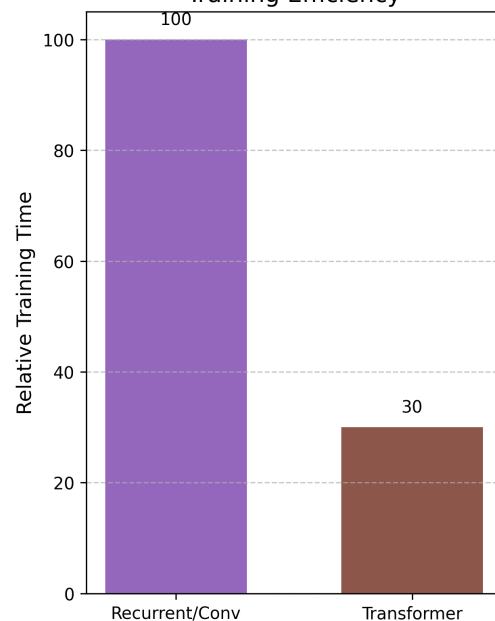


Transformer Performance vs. Traditional Models  
WMT 2014 English-to-German Translation



Training Efficiency



Scaled Dot-Product Attention Weights



Multi-Head Attention: Different Attention Patterns

