

Sequence Models: Practical Implementation Checklist

Essential considerations before deploying your model



Data Preparation

Tokenization

Choose appropriate tokenizer (BPE, WordPiece, SentencePiece)

Vocabulary Size

Balance between coverage and efficiency (typically 10K-50K)

Special Tokens

Define PAD, UNK, BOS, EOS tokens correctly

Sequence Length

Set max length based on data distribution (e.g., 128, 256, 512)



Model Architecture

Hidden Dimensions

Choose appropriate size (256, 512, 1024) based on task

Number of Layers

Balance depth and training time (2-12 layers typical)

Attention Mechanism

Implement proper attention with correct masking

Dropout Rates

Add dropout for regularization (0.1-0.3 typical)



Training Strategy

Teacher Forcing

Use during training, consider scheduled sampling

Batch Size

Optimize for GPU memory (32-128 typical)



Implementation Details

Padding & Masking

Properly handle variable-length sequences

Loss Calculation

Ignore padding tokens in loss computation

Learning Rate

Use warmup + decay schedule (e.g., 1e-4 to 1e-5)

 Gradient Clipping

Prevent exploding gradients (clip norm 1.0-5.0)

 Inference Strategy

Choose beam search, greedy, or sampling decoding

 Evaluation Metrics

Track BLEU, perplexity, or task-specific metrics

 **Quick Reference: Typical Hyperparameters****Hidden Size**

256-1024

**Batch Size**

32-128

**Vocab Size**

10K-50K

**Learning Rate**

1e-4 ~ 1e-3