

Chapter 5

Pretraining Large Language Models: Evaluation, Training, and Text Generation

CSC 375 - Generative AI

Prof. Rongyu Lin
Quinnipiac University
School of Computing and Engineering

Required Reading

Required Readings Before Class:

- **Raschka:** "Build a Large Language Model (From Scratch)" - Chapter 5, Pages 150-190
- **Focus Topics:** Loss calculation, training loops, text generation strategies, pretrained weights
- **Supplementary:** PyTorch documentation for cross-entropy loss and optimization algorithms

Before We Begin: A Question to Ponder

Imagine you have just built a sophisticated GPT model with millions of parameters, but when you ask it to complete "The weather today is...", it responds with random gibberish like "rentingetic wouldn? refres RexMeCHicular". Your model has the right architecture, the right size, and all the computational power needed—yet it produces meaningless text.

This scenario illustrates the fundamental challenge we face in Chapter 5: transforming a randomly initialized model into an intelligent language generator. How do we measure progress objectively? How do we systematically improve the model's predictions? And once trained, how do we control the creativity and coherence of the generated text?

Opening Reflection: Before reading further, consider this: How would you define "good" text generation mathematically? What would success look like in numerical terms? Think about the relationship between prediction accuracy and language understanding.

1 Learning Objectives

By the end of this lecture, you will be able to:

1. **Evaluate** generative models using cross-entropy loss and understand the mathematical foundation
2. **Implement** complete training loops for LLMs with proper data handling and optimization
3. **Calculate** and interpret training/validation losses to monitor model performance

4. **Apply** advanced text generation techniques including temperature and top-k sampling
5. **Load** and utilize pretrained model weights for transfer learning applications
6. **Understand** the complete pretraining process from initialization to deployment

2 The Training Pipeline: From Random to Remarkable

2.1 Understanding Weight Parameters in Large Language Models

Large Language Models derive their power from millions or billions of trainable parameters—the weights that determine how information flows through the neural network. Understanding these parameters is crucial for effective training and optimization.

What Are Weight Parameters? Weight parameters are numerical values that control the strength of connections between neurons in the neural network. During training, these weights are systematically adjusted to minimize prediction errors, gradually teaching the model to understand and generate coherent text.

Location of Parameters in LLMs:

- **Embedding Layers:** Convert token IDs to dense vector representations
- **Attention Mechanism Weights:** Query, key, value, and output projection matrices
- **Feed-Forward Network Parameters:** Two linear transformations with activation functions
- **Layer Normalization:** Scale and shift parameters for stabilizing training
- **Position Encoding:** Learned or fixed parameters for sequence order information

Accessing Parameters in PyTorch:

```

1 # Count total parameters in your model
2 total_params = sum(p.numel() for p in model.parameters())
3 print(f"Total parameters: {total_params:,}")
4
5 # Access specific layer parameters
6 for name, param in model.named_parameters():
7     print(f"{name}: {param.shape}")
8
9 # Example output for GPT-2 Small:
10 # transformer.wte.weight: torch.Size([50257, 768]) # Token embeddings
11 # transformer.wpe.weight: torch.Size([1024, 768]) # Position embeddings
12 # transformer.h.0.attn.c_attn.weight: torch.Size([768, 2304]) # Attention
    weights

```

Scale Perspectives (Reference: Raschka p. 151):

Model	Parameters	Memory (FP32)
GPT-2 Small	124 million	0.5 GB
GPT-2 Medium	350 million	1.4 GB
GPT-2 Large	774 million	3.1 GB
GPT-3	175 billion	700 GB

Scale Reality Check: The jump from GPT-2 to GPT-3 represents a 226x increase in parameters. This dramatic scaling, combined with more training data and compute, leads to the emergence of new capabilities that smaller models cannot achieve.

2.2 The Challenge: Untrained Model Behavior

Before training, even the most sophisticated language model produces incoherent text because its weights are randomly initialized. Understanding this baseline helps us appreciate the magnitude of the training challenge.

Example of Untrained Model Output: When prompted with "Every effort moves you", an untrained GPT model might generate:

"Every effort moves you rentingetic wouldn? refres RexMeCHicular stren"

Why Untrained Models Fail:

1. **Random Weight Initialization:** Parameters start with random values, typically drawn from normal distributions
2. **No Learned Patterns:** The model has never seen language data and knows nothing about word relationships
3. **Uniform Probability Distributions:** Without training, all tokens in the vocabulary have roughly equal probability
4. **No Contextual Understanding:** The attention mechanisms haven't learned to focus on relevant context

Testing Untrained Model Behavior:

```
1 # Generate text with untrained model
2 start_context = "Every effort moves you"
3 model.eval() # Set to evaluation mode
4 context = torch.tensor(encode_text(start_context)).unsqueeze(0)
5
6 # Use the simple generation function
7 generated = generate_text_simple(
8     model,
9     context,
10    max_new_tokens=10,
11    context_size=GPT_CONFIG_124M["context_length"]
12 )
13 print(f"Untrained output: {generated}")
```

The Training Imperative: The incoherent output from untrained models demonstrates why we need systematic training with proper loss functions. The goal is to transform random parameter values into meaningful language representations through exposure to vast amounts of text data.

3 Model Evaluation: Measuring Language Understanding

3.1 The Evaluation Challenge

Evaluating text generation quality presents unique challenges compared to traditional machine learning tasks. Unlike classification problems with clear right/wrong answers, language generation involves subjective qualities like coherence, creativity, and appropriateness.

Human Evaluation Limitations:

- **Inconsistency:** Different evaluators may rate the same text differently
- **Scalability:** Manual evaluation is time-consuming and expensive
- **Subjectivity:** Quality assessment varies based on context and personal preferences

- **Non-differentiable:** Human judgments cannot drive gradient-based optimization

The Solution: Cross-Entropy Loss Cross-entropy loss provides an objective, mathematical measure of model performance based on next-token prediction accuracy. A model that consistently assigns high probability to the actual next tokens in training text demonstrates better language understanding.

3.2 Cross-Entropy Loss: Mathematical Foundation

Cross-entropy loss measures the difference between predicted probability distributions and the true distribution of language. In language modeling, the "true" distribution assigns probability 1 to the actual next token and 0 to all others.

Mathematical Definition:

$$\text{CrossEntropy} = - \sum_{i=1}^V y_i \log(\hat{y}_i)$$

Where:

- V = vocabulary size (typically 50,000+ tokens)
- y_i = true probability (1 for correct token, 0 for all others)
- \hat{y}_i = predicted probability from the model

Simplified Form for Language Modeling: Since only one token has true probability 1, this simplifies to:

$$\text{CrossEntropy} = -\log(\hat{y}_{\text{target}})$$

Where \hat{y}_{target} is the model's predicted probability for the actual next token.

Intuitive Understanding:

- **Low Loss (Good):** High probability assigned to correct token
 - Example: $P(\text{"cat"} \rightarrow \text{"The"}) = 0.9 \rightarrow \text{Loss} = -\log(0.9) = 0.046$
- **High Loss (Bad):** Low probability assigned to correct token
 - Example: $P(\text{"cat"} \rightarrow \text{"The"}) = 0.1 \rightarrow \text{Loss} = -\log(0.1) = 1.0$

3.3 Step-by-Step Loss Calculation

Let's walk through a complete example of calculating cross-entropy loss for language modeling using real tensors and operations.

Example Setup:

```
1 # Two training sequences
2 inputs = torch.tensor([[16833, 3626, 6100], [40, 1107, 588]])
3 targets = torch.tensor([[3626, 6100, 345], [1107, 588, 11311]])
4
5 # Note: targets are inputs shifted by one position
6 # This creates the next-token prediction task
```

Step 1: Model Forward Pass

```
1 # Get logits from model
2 logits = model(inputs) # Shape: [batch_size, seq_len, vocab_size]
3 print(f"Logits shape: {logits.shape}") # [2, 3, 50257]
```

Step 2: Convert Logits to Probabilities Logits are raw model outputs that can be any real numbers. We use softmax to convert them to proper probabilities:

$$\text{Softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^V e^{z_j}}$$

```
1 # Apply softmax transformation
2 probas = torch.softmax(logits, dim=-1)
3 print(f"Probabilities shape: {probas.shape}")
4 print(f"Sum check: {probas[0, 0].sum()}") # Should be 1.0
5
6 # Properties of softmax output:
7 # - All values between 0 and 1
8 # - Sum of all probabilities = 1
9 # - Larger logits  $\rightarrow$  higher probabilities
```

Step 3: Extract Target Probabilities We need to extract the probabilities that the model assigned to the actual next tokens:

```
1 # Get probabilities for target tokens using advanced indexing
2 # This extracts the probability at each target token position
3 target_probas = probas[
4     range(len(targets)),          # Batch dimension
5     range(targets.shape[1]),      # Sequence dimension
6     targets.flatten()             # Use target IDs as indices
7 ]
8
9 # Reshape to match original structure
10 target_probas_resaped = target_probas.reshape(targets.shape)
11 print(f"Target probabilities: {target_probas_resaped}")
```

Step 4: Calculate Cross-Entropy Loss

```
1 # Apply negative logarithm and average
2 log_probas = torch.log(target_probas_resaped)
3 neg_log_probas = -log_probas
4 loss = neg_log_probas.mean()
5
6 print(f"Log probabilities: {log_probas}")
7 print(f"Negative log probabilities: {neg_log_probas}")
8 print(f"Cross-entropy loss: {loss.item()}")
```

Why Logarithms?

- **Mathematical Stability:** Prevents numerical underflow with very small probabilities
- **Computational Efficiency:** Converts products to sums for easier calculation
- **Optimization Properties:** Better gradient behavior during backpropagation
- **Information Theory:** Logarithms measure "surprise" or information content

3.4 PyTorch Cross-Entropy Function

While understanding the manual calculation is important, PyTorch provides an optimized implementation that handles edge cases and numerical stability:

```
1 import torch.nn.functional as F
2
3 # Manual calculation (what we just learned)
4 manual_loss = neg_log_probas.mean()
5
```

```

6 # PyTorch built-in function
7 # Note: Takes logits directly, not probabilities
8 pytorch_loss = F.cross_entropy(
9     logits.flatten(0, 1),      # Reshape to [batch*seq, vocab]
10    targets.flatten()          # Reshape to [batch*seq]
11 )
12
13 print(f"Manual calculation: {manual_loss:.6f}")
14 print(f"PyTorch function: {pytorch_loss:.6f}")
15 print(f"Difference: {abs(manual_loss - pytorch_loss):.10f}")

```

Advantages of PyTorch's Function:

- **Numerical Stability:** Uses log-sum-exp trick to prevent overflow/underflow
- **Efficiency:** Optimized C++ implementation with CUDA support
- **Convenience:** Takes logits directly, no need for manual softmax
- **Edge Case Handling:** Properly handles zero probabilities and extreme values
- **Memory Efficiency:** Avoids creating intermediate probability tensors

Best Practice: Always use PyTorch's built-in functions for training, but understand the underlying mathematics for debugging and theoretical comprehension.

4 Training Data Preparation and Management

4.1 Training vs. Validation Split Strategy

Proper data splitting is crucial for monitoring model generalization and preventing overfitting. The validation set serves as an independent measure of model performance throughout training.

Split Rationale (Reference: Raschka p. 162):

- **Training Set (90%):** Used for gradient calculation and weight updates
- **Validation Set (10%):** Used only for performance monitoring, never for training
- **Temporal Ordering:** Validation data typically comes from later time periods to simulate real-world deployment

Implementation:

```

1 # Split the text data
2 total_size = len(text_data)
3 train_ratio = 0.90
4 split_idx = int(train_ratio * total_size)
5
6 train_data = text_data[:split_idx]
7 val_data = text_data[split_idx:]
8
9 print(f"Training characters: {len(train_data):,}")
10 print(f"Validation characters: {len(val_data):,}")
11 print(f"Split ratio: {len(train_data)/len(text_data):.2%}")
12
13 # Verify the split preserves data integrity
14 assert len(train_data) + len(val_data) == len(text_data)
15 assert train_data[-10:] + val_data[:10] == text_data[split_idx-10:split_idx+10]

```

Monitoring Generalization:

- **Good Generalization:** Validation loss decreases alongside training loss
- **Overfitting:** Training loss continues decreasing while validation loss increases
- **Underfitting:** Both losses remain high and plateau early

4.2 Efficient Data Loader Implementation

Data loaders handle the complex task of converting raw text into batched tensors suitable for training. Proper implementation is crucial for training efficiency and memory management.

Key Parameters:

- **batch_size:** Number of sequences processed simultaneously (affects memory usage)
- **max_length:** Maximum sequence length (context window)
- **stride:** Overlap between consecutive sequences (prevents context loss)
- **shuffle:** Randomize order for training (True), preserve order for validation (False)

Data Loader Creation:

```

1 # Configuration for training
2 batch_size = 2      # Small for learning purposes, typically 16-128 in practice
3 max_length = 256    # Context window size
4 stride = 128        # 50% overlap between sequences
5
6 # Create training data loader
7 train_loader = create_dataloader_v1(
8     train_data,
9     batch_size=batch_size,
10    max_length=max_length,
11    stride=stride,
12    shuffle=True     # Randomize for better generalization
13 )
14
15 # Create validation data loader
16 val_loader = create_dataloader_v1(
17     val_data,
18     batch_size=batch_size,
19     max_length=max_length,
20     stride=stride,
21     shuffle=False   # Preserve order for consistent evaluation
22 )
23
24 # Inspect data loader output
25 for batch in train_loader:
26     print(f"Batch shape: {batch.shape}") # [batch_size, max_length]
27     break

```

Memory and Performance Considerations:

- **Batch Size Trade-offs:** Larger batches use more memory but provide more stable gradients
- **Sequence Length:** Longer sequences capture more context but require quadratically more attention computation
- **Stride Strategy:** Smaller strides increase data coverage but create more computational overhead

5 Training Loop Implementation

5.1 Complete Training Loop Architecture

The training loop is the heart of the learning process, orchestrating forward passes, loss calculation, backpropagation, and weight updates in a systematic cycle.

Training Loop Components:

1. **Data Loading:** Fetch batches from training data loader
2. **Forward Pass:** Compute model predictions and loss
3. **Backward Pass:** Calculate gradients via backpropagation
4. **Optimization:** Update model weights using computed gradients
5. **Monitoring:** Track performance on validation data
6. **Logging:** Record metrics for analysis and visualization

Training Loop Structure:

```
1 def train_model_simple(model, train_loader, val_loader, optimizer, device,
2                        num_epochs, eval_freq, eval_iter, start_context):
3     """
4     Complete training loop with monitoring and evaluation
5
6     Args:
7         model: The language model to train
8         train_loader: DataLoader for training data
9         val_loader: DataLoader for validation data
10        optimizer: Optimization algorithm (e.g., AdamW)
11        device: Computing device (CPU/GPU)
12        num_epochs: Number of training epochs
13        eval_freq: How often to evaluate on validation set
14        eval_iter: Number of batches for validation evaluation
15        start_context: Text prompt for generation during training
16    """
17
18    # Tracking variables
19    train_losses, val_losses, track_tokens_seen = [], [], []
20    tokens_seen, global_step = 0, -1
21
22    # Training loop
23    for epoch in range(num_epochs):
24        model.train() # Set to training mode
25
26        for input_batch, target_batch in train_loader:
27            optimizer.zero_grad() # Clear previous gradients
28
29            # Move data to device
30            input_batch = input_batch.to(device)
31            target_batch = target_batch.to(device)
32
33            # Forward pass
34            logits = model(input_batch)
35            loss = torch.nn.functional.cross_entropy(
36                logits.flatten(0, 1),
37                target_batch.flatten()
38            )
39
40            # Backward pass
41            loss.backward() # Compute gradients
```



```

42         optimizer.step() # Update weights
43
44     # Update tracking
45     tokens_seen += input_batch.numel()
46     global_step += 1
47
48     # Periodic evaluation
49     if global_step % eval_freq == 0:
50         train_loss, val_loss = evaluate_model(
51             model, train_loader, val_loader, device, eval_iter
52         )
53         train_losses.append(train_loss)
54         val_losses.append(val_loss)
55         track_tokens_seen.append(tokens_seen)
56
57         print(f"Ep {epoch+1} (Step {global_step:06d}): "
58               f"Train loss {train_loss:.3f}, Val loss {val_loss:.3f}")
59
60     return train_losses, val_losses, track_tokens_seen

```

5.2 Optimizer Selection: AdamW

The choice of optimizer significantly impacts training dynamics and final model performance. AdamW (Adam with Weight Decay) is the standard choice for transformer model training.

AdamW Advantages:

- **Adaptive Learning Rates:** Different learning rates for each parameter based on gradient history
- **Momentum:** Accumulates gradient direction over time for smoother optimization
- **Weight Decay:** Proper L2 regularization to prevent overfitting
- **Stability:** Robust performance across different model sizes and datasets

AdamW Setup:

```

1 import torch.optim as optim
2
3 # Configure AdamW optimizer
4 optimizer = optim.AdamW(
5     model.parameters(),
6     lr=0.0004,           # Learning rate (typical range: 1e-5 to 1e-3)
7     weight_decay=0.1     # L2 regularization strength
8 )
9
10 # Learning rate scheduling (optional but recommended)
11 scheduler = optim.lr_scheduler.CosineAnnealingLR(
12     optimizer,
13     T_max=num_epochs,    # Period of cosine decay
14     eta_min=1e-6         # Minimum learning rate
15 )

```

Key Hyperparameters (Reference: Raschka p. 166):

Parameter	Typical Range	Our Choice
Learning Rate	1e-5 to 1e-3	4e-4
Weight Decay	0.01 to 0.3	0.1
Beta1	0.9 to 0.95	0.9
Beta2	0.95 to 0.999	0.999

5.3 Training Progress Monitoring

Effective monitoring during training provides insights into model behavior and helps identify issues early in the training process.

Evaluation Function:

```
1 def evaluate_model(model, train_loader, val_loader, device, eval_iter):
2     """
3     Evaluate model performance on training and validation sets
4
5     Args:
6         model: Model to evaluate
7         train_loader: Training data loader
8         val_loader: Validation data loader
9         device: Computing device
10        eval_iter: Number of batches to evaluate
11
12    Returns:
13        train_loss: Average loss on training set
14        val_loss: Average loss on validation set
15    """
16    model.eval() # Set to evaluation mode
17
18    with torch.no_grad(): # Disable gradient computation
19        train_loss = calc_loss_loader(train_loader, model, device,
20                                     num_batches=eval_iter)
21        val_loss = calc_loss_loader(val_loader, model, device,
22                                   num_batches=eval_iter)
23
24    model.train() # Return to training mode
25    return train_loss, val_loss
26
27 def calc_loss_loader(data_loader, model, device, num_batches=None):
28     """Calculate average loss over specified number of batches"""
29     total_loss = 0.0
30     if num_batches is None:
31         num_batches = len(data_loader)
32     else:
33         num_batches = min(num_batches, len(data_loader))
34
35     for i, (input_batch, target_batch) in enumerate(data_loader):
36         if i >= num_batches:
37             break
38
39         input_batch = input_batch.to(device)
40         target_batch = target_batch.to(device)
41
42         logits = model(input_batch)
43         loss = torch.nn.functional.cross_entropy(
44             logits.flatten(0, 1),
45             target_batch.flatten()
46         )
47         total_loss += loss.item()
48
49     return total_loss / num_batches
```

Key Monitoring Metrics:

- **Training Loss:** Primary optimization target (should consistently decrease)
- **Validation Loss:** Generalization indicator (should track training loss)
- **Generation Quality:** Periodic text samples to assess coherence

- **Training Speed:** Tokens processed per second for efficiency tracking

6 Advanced Text Generation Strategies

6.1 Beyond Deterministic Generation

Simple text generation uses greedy decoding—always selecting the token with highest probability. While deterministic and fast, this approach often produces repetitive and boring text. Advanced sampling strategies introduce controlled randomness for more diverse and engaging outputs.

Limitations of Greedy Decoding:

- **Repetition:** Tends to produce repetitive phrases and ideas
- **Lack of Creativity:** Always chooses the "safest" option
- **Exposure Bias:** Training uses teacher forcing, but generation is autoregressive
- **Local Optima:** May get stuck in suboptimal continuation paths

6.2 Temperature Sampling: Controlling Creativity

Temperature sampling modifies the probability distribution before sampling, allowing fine-grained control over generation randomness and creativity.

Mathematical Foundation: Temperature scaling adjusts logits before applying softmax:

$$\text{Softmax}_T(z_i) = \frac{e^{z_i/T}}{\sum_{j=1}^V e^{z_j/T}}$$

Where T is the temperature parameter controlling the "sharpness" of the distribution.

Temperature Effects:

- **$T < 1$ (Conservative):** Sharpens distribution, more deterministic
 - $T = 0.1$: Very focused, minimal creativity
 - $T = 0.5$: Somewhat focused, reduced randomness
- **$T = 1$ (Default):** Original model distribution, balanced
- **$T > 1$ (Creative):** Flattens distribution, more random
 - $T = 1.5$: Increased creativity, some randomness
 - $T = 2.0$: High creativity, significant randomness

Temperature Sampling Implementation:

```

1 def generate_with_temperature(model, tokenizer, prompt, max_tokens=50,
2                               temperature=1.0, context_size=1024):
3     """
4     Generate text using temperature sampling
5
6     Args:
7         model: Trained language model
8         tokenizer: Text tokenizer
9         prompt: Starting text
10        max_tokens: Maximum tokens to generate
11        temperature: Sampling temperature
12        context_size: Maximum context length

```

```

13
14     Returns:
15         Generated text string
16     """
17     model.eval()
18
19     # Convert prompt to tokens
20     context = torch.tensor(tokenizer.encode(prompt)).unsqueeze(0)
21
22     with torch.no_grad():
23         for _ in range(max_tokens):
24             # Truncate context if too long
25             context_trunc = context[:, -context_size:]
26
27             # Get logits from model
28             logits = model(context_trunc)
29
30             # Apply temperature scaling
31             logits = logits[:, -1, :] / temperature
32
33             # Convert to probabilities
34             probabilities = torch.softmax(logits, dim=-1)
35
36             # Sample next token
37             next_token = torch.multinomial(probabilities, num_samples=1)
38
39             # Append to context
40             context = torch.cat([context, next_token], dim=1)
41
42     # Convert back to text
43     return tokenizer.decode(context.squeeze(0).tolist())

```

Temperature Selection Guidelines:

Use Case	Temperature	Characteristics
Technical Writing	0.1 - 0.3	Focused, accurate, minimal creativity
News Articles	0.3 - 0.7	Balanced accuracy and readability
Creative Writing	0.7 - 1.2	Good balance of coherence and creativity
Poetry/Art	1.0 - 2.0	High creativity, experimental language

6.3 Top-k Sampling: Focused Creativity

Top-k sampling combines the benefits of probabilistic sampling with quality control by restricting sampling to the k most likely tokens.

Algorithm:

1. Compute token probabilities from model logits
2. Sort tokens by probability (highest to lowest)
3. Keep only the top-k most likely tokens
4. Set probabilities of other tokens to zero
5. Renormalize remaining probabilities to sum to 1
6. Sample from this restricted distribution

Top-k Implementation:

```
1 def sample_top_k(logits, k=10, temperature=1.0):
2     """
3     Sample from top-k most likely tokens
4
5     Args:
6         logits: Model output logits [vocab_size]
7         k: Number of top tokens to consider
8         temperature: Sampling temperature
9
10    Returns:
11        Sampled token ID
12    """
13    # Apply temperature scaling
14    logits = logits / temperature
15
16    # Get top-k logits and indices
17    top_logits, top_indices = torch.topk(logits, k)
18
19    # Convert to probabilities
20    top_probabilities = torch.softmax(top_logits, dim=-1)
21
22    # Sample from top-k
23    sample_idx = torch.multinomial(top_probabilities, num_samples=1)
24
25    # Return original token index
26    return top_indices[sample_idx]
```

Optimal k Values:

- **k = 1:** Equivalent to greedy decoding (deterministic)
- **k = 10-20:** Good balance for most applications
- **k = 50-100:** More creative, suitable for creative writing
- **k = vocab_size:** Equivalent to pure temperature sampling

6.4 Combined Sampling Strategy

The most effective approach often combines temperature and top-k sampling to achieve both quality control and creative flexibility.

Advanced Generation Function:

```
1 def generate_advanced(model, tokenizer, prompt, max_tokens=100,
2                       temperature=1.0, top_k=50, context_size=1024):
3     """
4     Advanced text generation with temperature and top-k sampling
5
6     Args:
7         model: Trained language model
8         tokenizer: Text tokenizer
9         prompt: Starting prompt text
10        max_tokens: Maximum number of tokens to generate
11        temperature: Sampling temperature (controls randomness)
12        top_k: Number of top tokens to consider
13        context_size: Maximum context window size
14
15    Returns:
16        Generated text as string
17    """
18    model.eval()
```

```

19
20 # Tokenize initial prompt
21 context = torch.tensor(tokenizer.encode(prompt)).unsqueeze(0)
22
23 generated_tokens = []
24
25 with torch.no_grad():
26     for _ in range(max_tokens):
27         # Manage context window
28         context_input = context[:, -context_size:]
29
30         # Forward pass
31         logits = model(context_input)
32         logits = logits[:, -1, :] # Get last position logits
33
34         # Apply temperature
35         if temperature != 1.0:
36             logits = logits / temperature
37
38         # Apply top-k filtering
39         if top_k > 0:
40             # Get top-k values and indices
41             top_logits, top_indices = torch.topk(logits,
42                                                 min(top_k, logits.size(-1)))
43
44             # Create filtered logits tensor
45             filtered_logits = torch.full_like(logits, float('-inf'))
46             filtered_logits.scatter_(1, top_indices, top_logits)
47             logits = filtered_logits
48
49         # Convert to probabilities and sample
50         probabilities = torch.softmax(logits, dim=-1)
51         next_token = torch.multinomial(probabilities, num_samples=1)
52
53         # Add to context and generated tokens
54         context = torch.cat([context, next_token], dim=1)
55         generated_tokens.append(next_token.item())
56
57         # Stop if end-of-text token generated
58         if next_token.item() == tokenizer.encode('<|endoftext|>')[0]:
59             break
60
61 # Decode generated tokens
62 generated_text = tokenizer.decode(generated_tokens)
63 return prompt + generated_text

```

Sampling Strategy Guidelines:

Application	Temperature	Top-k	Expected Quality
Code Generation	0.1	10	High precision, low creativity
Technical Writing	0.3	20	Accurate, professional
General Chat	0.7	40	Balanced, engaging
Creative Writing	1.0	60	Creative, diverse
Experimental Art	1.5	100	Highly creative, unpredictable

7 Loading and Using Pretrained Weights

7.1 The Pretraining Advantage

Training large language models from scratch requires enormous computational resources—millions of dollars and months of training time. Pretrained models offer a practical alternative by pro-

viding weights learned from massive datasets that can be fine-tuned for specific applications.

Benefits of Pretrained Models:

- **Cost Efficiency:** Avoid expensive pretraining computation
- **Time Savings:** Start with learned language representations
- **Performance:** Benefit from training on massive, diverse datasets
- **Accessibility:** Makes advanced AI capabilities available to smaller organizations

Transfer Learning Process:

1. **Pretraining:** Large-scale training on diverse text (e.g., web crawl data)
2. **Fine-tuning:** Adapt pretrained model to specific tasks or domains
3. **Evaluation:** Test performance on target applications
4. **Deployment:** Use fine-tuned model in production systems

7.2 OpenAI GPT-2 Weights

OpenAI provides pretrained GPT-2 models in multiple sizes, offering a balance between performance and computational requirements.

Available Model Sizes:

Model	Parameters	Layers	Hidden Size
GPT-2 117M	117 million	12	768
GPT-2 345M	345 million	24	1024
GPT-2 762M	762 million	36	1280
GPT-2 1542M	1.5 billion	48	1600

Model Configuration Alignment: Before loading pretrained weights, ensure your model configuration exactly matches the pretrained model architecture:

```
1 # Configuration for GPT-2 117M (small)
2 GPT_CONFIG_124M = {
3     "vocab_size": 50257,      # Vocabulary size
4     "context_length": 1024,   # Maximum sequence length
5     "emb_dim": 768,           # Embedding dimension
6     "n_heads": 12,            # Number of attention heads
7     "n_layers": 12,           # Number of transformer layers
8     "drop_rate": 0.1,         # Dropout rate
9     "qkv_bias": False         # Query-key-value bias
10 }
11
12 # Verify configuration matches pretrained model
13 def verify_config(config, model_name="gpt2"):
14     """Verify configuration matches pretrained model requirements"""
15     if model_name == "gpt2": # 117M model
16         assert config["vocab_size"] == 50257
17         assert config["context_length"] == 1024
18         assert config["emb_dim"] == 768
19         assert config["n_heads"] == 12
20         assert config["n_layers"] == 12
21         print("\checkmark\ Configuration verified for GPT-2 117M")
22     else:
23         raise ValueError(f"Unknown model: {model_name}")
24
25 verify_config(GPT_CONFIG_124M)
```

7.3 Weight Loading Process

Loading pretrained weights requires careful attention to tensor shapes, parameter names, and model architecture compatibility.

Weight Loading Implementation:

```
1 import requests
2 import json
3 import numpy as np
4 from tqdm import tqdm
5
6 def download_and_load_gpt2(model_size="124M", models_dir="gpt2"):
7     """
8     Download and load GPT-2 pretrained weights
9
10    Args:
11        model_size: Model size ("124M", "355M", "774M", "1558M")
12        models_dir: Directory to store model files
13
14    Returns:
15        Dictionary containing model weights and configuration
16    """
17
18    # Create model directory
19    os.makedirs(models_dir, exist_ok=True)
20
21    # Download model files
22    base_url = "https://openaipublic.blob.core.windows.net/gpt-2/models"
23
24    for filename in ["checkpoint", "encoder.json", "hparams.json",
25                    "model.ckpt.data-00000-of-00001", "model.ckpt.index",
26                    "model.ckpt.meta", "vocab.bpe"]:
27
28        file_url = f"{base_url}/{model_size}/{filename}"
29        file_path = os.path.join(models_dir, filename)
30
31        if not os.path.exists(file_path):
32            print(f"Downloading {filename}...")
33            response = requests.get(file_url)
34            with open(file_path, "wb") as f:
35                f.write(response.content)
36
37    # Load weights from TensorFlow checkpoint
38    import tensorflow as tf
39
40    tf_ckpt_path = os.path.join(models_dir, "model.ckpt")
41    tf_weights = {}
42
43    for name, _ in tf.train.list_variables(tf_ckpt_path):
44        array = tf.train.load_variable(tf_ckpt_path, name)
45        tf_weights[name] = array
46
47    return tf_weights
48
49 def load_weights_into_gpt(gpt_model, weights_dict):
50     """
51     Load pretrained weights into GPT model
52
53    Args:
54        gpt_model: PyTorch GPT model instance
55        weights_dict: Dictionary of pretrained weights
56    """
57
```



```

58 # Weight mapping from TensorFlow to PyTorch naming
59 weight_mapping = {
60     "wte": "transformer.wte.weight",           # Token embeddings
61     "wpe": "transformer.wpe.weight",           # Position embeddings
62     "ln_f/g": "transformer.ln_f.weight",        # Final layer norm gamma
63     "ln_f/b": "transformer.ln_f.bias",         # Final layer norm beta
64 }
65
66 # Load layer-specific weights
67 for layer_idx in range(gpt_model.config["n_layers"]):
68     # Attention weights
69     weight_mapping.update({
70         f"h{layer_idx}/attn/c_attn/w": f"transformer.h.{layer_idx}.attn.
71             c_attn.weight",
72         f"h{layer_idx}/attn/c_attn/b": f"transformer.h.{layer_idx}.attn.
73             c_attn.bias",
74         f"h{layer_idx}/attn/c_proj/w": f"transformer.h.{layer_idx}.attn.
75             c_proj.weight",
76         f"h{layer_idx}/attn/c_proj/b": f"transformer.h.{layer_idx}.attn.
77             c_proj.bias",
78
79         # Feed-forward weights
80         f"h{layer_idx}/mlp/c_fc/w": f"transformer.h.{layer_idx}.mlp.c_fc.
81             weight",
82         f"h{layer_idx}/mlp/c_fc/b": f"transformer.h.{layer_idx}.mlp.c_fc.
83             bias",
84         f"h{layer_idx}/mlp/c_proj/w": f"transformer.h.{layer_idx}.mlp.
85             c_proj.weight",
86         f"h{layer_idx}/mlp/c_proj/b": f"transformer.h.{layer_idx}.mlp.
87             c_proj.bias",
88
89         # Layer normalization
90         f"h{layer_idx}/ln_1/g": f"transformer.h.{layer_idx}.ln_1.weight",
91         f"h{layer_idx}/ln_1/b": f"transformer.h.{layer_idx}.ln_1.bias",
92         f"h{layer_idx}/ln_2/g": f"transformer.h.{layer_idx}.ln_2.weight",
93         f"h{layer_idx}/ln_2/b": f"transformer.h.{layer_idx}.ln_2.bias",
94     })
95
96 # Load weights into model
97 with torch.no_grad():
98     for tf_name, pytorch_name in weight_mapping.items():
99         if tf_name in weights_dict:
100             tf_tensor = torch.from_numpy(weights_dict[tf_name])
101
102             # Handle weight transpose for linear layers
103             if "c_attn/w" in tf_name or "c_proj/w" in tf_name or "c_fc/w"
104                 in tf_name:
105                 tf_tensor = tf_tensor.t() # Transpose for PyTorch format
106
107             # Get corresponding PyTorch parameter
108             pytorch_param = gpt_model.state_dict()[pytorch_name]
109
110             # Verify shapes match
111             assert tf_tensor.shape == pytorch_param.shape, \
112                 f"Shape mismatch for {pytorch_name}: {tf_tensor.shape} vs {
113                     pytorch_param.shape}"
114
115             # Copy weights
116             pytorch_param.copy_(tf_tensor)
117
118 print("\checkmark\ Pretrained weights loaded successfully")

```

7.4 Model Size and Performance Comparison

Understanding the trade-offs between different model sizes helps in selecting the appropriate pretrained model for specific applications.

Performance vs. Resource Trade-offs:

Model	Memory (GB)	Inference Speed	Quality	Use Case
GPT-2 117M	0.5	Very Fast	Good	Development, testing
GPT-2 345M	1.4	Fast	Better	Production apps
GPT-2 762M	3.1	Medium	Very Good	High-quality applications
GPT-2 1542M	6.2	Slower	Excellent	Research, premium apps

Model Selection Guidelines:

- **Prototyping:** Start with 117M for fast iteration
- **Production:** Use 345M for most commercial applications
- **High Quality:** Choose 762M for demanding tasks
- **Research:** Use 1542M for maximum performance

7.5 Pretrained Model Capabilities

Pretrained GPT-2 models demonstrate remarkable capabilities across various language tasks without task-specific fine-tuning.

Zero-Shot Capabilities:

```
1 # Test pretrained model capabilities
2 test_prompts = [
3     "Translate to French: Hello, how are you?",
4     "Q: What is the capital of France? A:",
5     "Write a Python function to calculate factorial:",
6     "Complete the sentence: The weather today is",
7     "Summarize: Large language models are AI systems..."
8 ]
9
10 def test_model_capabilities(model, tokenizer, prompts, max_tokens=50):
11     """Test various capabilities of pretrained model"""
12
13     model.eval()
14     results = {}
15
16     for prompt in prompts:
17         generated = generate_advanced(
18             model=model,
19             tokenizer=tokenizer,
20             prompt=prompt,
21             max_tokens=max_tokens,
22             temperature=0.7,
23             top_k=40
24         )
25
26         results[prompt] = generated
27         print(f"Prompt: {prompt}")
28         print(f"Generated: {generated}")
29         print("-" * 80)
30
31     return results
32
```

```
33 # Run capability tests
34 capabilities = test_model_capabilities(model, tokenizer, test_prompts)
```

Observed Capabilities:

- **Language Translation:** Basic translation between common languages
- **Question Answering:** Factual responses to straightforward questions
- **Code Generation:** Simple programming tasks and code completion
- **Text Completion:** Coherent continuation of various text types
- **Summarization:** Condensing longer texts into key points

8 Practical Applications and Deployment Considerations

8.1 Fine-tuning vs. Pretraining Decisions

Understanding when to use pretrained models versus training from scratch is crucial for project success and resource management.

Use Pretrained Models When:

- **Limited Computational Resources:** Budget constraints on training infrastructure
- **General Language Tasks:** Applications requiring broad language understanding
- **Fast Development Cycles:** Rapid prototyping and iteration needs
- **Standard Domains:** Working with common text types (news, conversations, etc.)

Consider Training from Scratch When:

- **Specialized Domains:** Highly technical or niche vocabularies
- **Privacy Requirements:** Sensitive data that cannot use external models
- **Specific Architectures:** Novel model designs not available pretrained
- **Ultimate Performance:** Maximum quality requirements with unlimited resources

Hybrid Approaches:

- **Domain Adaptation:** Continue pretraining on domain-specific data
- **Task-Specific Fine-tuning:** Adapt pretrained models to specific applications
- **Progressive Training:** Start with pretrained weights, then fine-tune with custom data

8.2 Ethical Considerations and Responsible Deployment

Deploying language models responsibly requires understanding their limitations and implementing appropriate safeguards.

Key Ethical Concerns:

- **Bias and Fairness:** Models reflect biases present in training data
- **Misinformation:** Potential for generating false or misleading content
- **Privacy:** Risk of reproducing sensitive information from training data

- **Transparency:** Users should understand they're interacting with AI systems

Mitigation Strategies:

- **Content Filtering:** Implement post-generation filtering for harmful content
- **User Education:** Clearly communicate AI limitations and appropriate usage
- **Monitoring Systems:** Track model outputs for quality and safety issues
- **Regular Evaluation:** Continuously assess model performance across diverse scenarios

9 Chapter 5 Summary and Key Takeaways

Essential Insights from Chapter 5:

1. **Objective Evaluation:** Cross-entropy loss provides a mathematical foundation for measuring and improving language model performance, bridging the gap between subjective text quality and quantitative optimization.
2. **Training Fundamentals:** Successful language model training requires careful orchestration of data preparation, loss calculation, optimization, and monitoring—each component crucial for learning language patterns.
3. **Advanced Generation:** Simple greedy decoding produces repetitive text; temperature and top-k sampling enable controllable creativity while maintaining coherence and quality.
4. **Pretrained Power:** Leveraging pretrained weights democratizes access to sophisticated language capabilities, enabling smaller organizations to build powerful applications without massive computational investments.
5. **Practical Deployment:** Real-world applications require balancing model capabilities with computational constraints, ethical considerations, and user expectations.

Technical Mastery Achieved:

- Understanding cross-entropy loss calculation and its role in optimization
- Implementing complete training loops with proper monitoring and evaluation
- Applying advanced sampling strategies for controlled text generation
- Loading and utilizing pretrained model weights effectively
- Making informed decisions about model selection and deployment strategies

Conceptual Understanding Developed:

- The relationship between mathematical optimization and language understanding
- How training transforms random parameters into meaningful representations
- The trade-offs between deterministic and probabilistic generation approaches
- The value proposition and limitations of transfer learning in NLP
- Ethical responsibilities in deploying language generation systems

10 Exercises and Practice Problems

Exercise 1: Loss Calculation Analysis Implement manual cross-entropy loss calculation for the following scenarios and compare with PyTorch's built-in function:

1. Calculate loss for a model that assigns uniform probability to all tokens
2. Calculate loss for a model with perfect next-token prediction
3. Analyze how loss changes as model predictions improve during training
4. Plot the relationship between prediction confidence and resulting loss values

Exercise 2: Training Loop Implementation Build a complete training loop with the following enhancements:

1. Add learning rate scheduling (cosine annealing or step decay)
2. Implement gradient clipping to prevent exploding gradients
3. Add early stopping based on validation loss plateaus
4. Include periodic model checkpointing for recovery
5. Implement training resume functionality from saved checkpoints

Exercise 3: Generation Strategy Comparison Compare different text generation strategies:

1. Generate 10 completions for the same prompt using different temperature values (0.1, 0.5, 1.0, 1.5, 2.0)
2. Test top-k sampling with k values of 1, 10, 50, and 100
3. Combine temperature and top-k sampling and evaluate the results
4. Design evaluation criteria for assessing generation quality
5. Analyze trade-offs between creativity and coherence

Exercise 4: Pretrained Model Analysis Explore pretrained model capabilities:

1. Load different sizes of GPT-2 (117M, 345M) and compare performance
2. Test zero-shot capabilities across various tasks
3. Measure inference speed and memory usage for different model sizes
4. Analyze the relationship between model size and output quality
5. Design benchmark tests for evaluating pretrained model performance

Exercise 5: Ethical AI Implementation Design responsible deployment strategies:

1. Identify potential biases in model outputs across different demographic groups
2. Implement content filtering systems for inappropriate or harmful text
3. Design user interfaces that clearly communicate AI limitations
4. Create monitoring systems for tracking model performance in production

5. Develop guidelines for appropriate use cases and applications

Next Class Preview: We'll dive deeper into the transformer architecture itself, building the attention mechanism from scratch and understanding how self-attention enables the powerful language modeling capabilities we've been utilizing.

Preparation for Lecture 6:

- **Raschka:** Chapter 3, Sections 3.1-3.5 (Pages 51-89) - Attention mechanisms
- **Review:** Linear algebra concepts (matrix multiplication, dot products)
- **Supplementary:** "Attention Is All You Need" paper introduction and overview