

How Loss is Calculated in Caption Generation

TL;DR

It IS classification! At each position, the model classifies which word (out of 4,374 options) comes next.

The Key Insight: Word-Level Classification

Your caption generation task:

Audio: [engine_sound.wav]

Target caption: "an engine is running loudly"

What the model actually does:

Position 1: Classify which word comes first

Options: ["a", "an", "the", "one", ...] (4,374 words)

Correct: "an"

Position 2: Classify which word comes second

Options: ["a", "an", "the", "engine", ...] (4,374 words)

Correct: "engine"

Position 3: Classify which word comes third

Options: ["is", "are", "was", ...] (4,374 words)

Correct: "is"

... and so on for each position

Each position = separate 4,374-way classification problem!

Step-by-Step Loss Calculation

Training Example

Input:

```
python
```

```
mel = [audio_features] # Audio: engine sound
caption = ["<sos>", "an", "engine", "is", "running", "loudly", "<eos>"]
```

Model forward pass:

```
python
```

```
logits = model(mel, caption)
# logits shape: [batch_size, seq_len, vocab_size]
# Example: [32, 7, 4374]
# 32 = batch size
# 7 = sequence length (7 words)
# 4374 = vocabulary size (number of possible words)
```

What are logits?

Logits = raw scores for each word at each position

```
python
```

```
# At position 0 (predicting "an"):
logits[0, 0, :] = [
    0.5,  # score for word_id 0 (<pad>)
    -1.2, # score for word_id 1 (<unk>)
    5.8,  # score for word_id 2 (<sos>)
    8.3,  # score for word_id 3 ("an") ← HIGH!
    2.1,  # score for word_id 4 ("a")
    ...
    0.3   # score for word_id 4373 ("zebra")
]

# Convert to probabilities with softmax:
probs = softmax(logits[0, 0, :])
# = [0.001, 0.0002, 0.02, 0.81, 0.05, ..., 0.0001]
# Model is 81% confident next word is "an" ✓
```

Teacher Forcing During Training

Important concept:

```
python
```

```
# During training, we use GROUND TRUTH as input
Input sequence: ["<sos>", "an", "engine", "is", "running"]
Target sequence: [    "an", "engine", "is", "running", "loudly"]
                  predict predict predict predict predict predict

# This is called "teacher forcing"
# Even if model predicts wrong word, we feed correct word next step
```

Why? Faster training, more stable learning

Cross-Entropy Loss Calculation

```
python
```

```
criterion = nn.CrossEntropyLoss(ignore_index=vocab['<pad>'])

# Reshape for loss calculation
logits_flat = logits.view(-1, vocab_size) # [batch*seq_len, vocab_size]
targets_flat = targets.view(-1)           # [batch*seq_len]

# Example:
# logits_flat[0] = scores for predicting word at position 0
# targets_flat[0] = 3 (correct word is "an" with id=3)

loss = criterion(logits_flat, targets_flat)
```

Mathematical Formula

For each position t :

$$\text{Loss}_t = -\log P(\text{correct_word}_t | \text{context}_t)$$

Where:

$$P(\text{word}_i) = \frac{e^{\text{logit}_i}}{\sum_{j=1}^V e^{\text{logit}_j}}$$

In plain English:

- Model outputs probability distribution over all 4,374 words
 - We want high probability on correct word
 - Loss = negative log probability of correct word
 - Lower loss = higher confidence in correct word
-

Concrete Example

Single word prediction:

```
python

# Ground truth: next word is "engine" (id=145)
target = 145

# Model outputs (simplified to 5 words for clarity):
logits = [1.2, # id=0: "a"
          8.5, # id=145: "engine" ← correct
          2.3, # id=200: "car"
          1.8, # id=300: "sound"
          0.5] # id=400: "music"

# Convert to probabilities
probs = softmax(logits) = [0.008, 0.977, 0.024, 0.015, 0.004]
                           ↑
                           97.7% confident!

# Cross-entropy loss
loss = -log(0.977) = 0.023 ← LOW (good prediction)

# If model was confused:
bad_probs = [0.3, 0.2, 0.3, 0.15, 0.05]
              ↑
              only 20% confident

bad_loss = -log(0.2) = 1.39 ← HIGH (bad prediction)
```

Full sequence:

```
python
```

Caption: "an engine is running"

Targets: [3, 145, 200, 450] # Word IDs

Position 0: Predict "an" (id=3)

Model confidence: 81%

$$\text{Loss}_0 = -\log(0.81) = 0.21$$

Position 1: Predict "engine" (id=145)

Model confidence: 65%

$$\text{Loss}_1 = -\log(0.65) = 0.43$$

Position 2: Predict "is" (id=200)

Model confidence: 92%

$$\text{Loss}_2 = -\log(0.92) = 0.08$$

Position 3: Predict "running" (id=450)

Model confidence: 73%

$$\text{Loss}_3 = -\log(0.73) = 0.31$$

Average loss across all positions

$$\text{Total Loss} = (0.21 + 0.43 + 0.08 + 0.31) / 4 = 0.26$$

Training vs Validation Loss

Training Loss

```
python
```

```
for mel, captions in train_loader:  
    # Forward pass  
    logits = model(mel, captions[:, :-1]) # Input: all but last word  
  
    # Compute loss against targets  
    targets = captions[:, 1:] # Target: all but first word  
  
    loss = criterion(  
        logits.reshape(-1, vocab_size),  
        targets.reshape(-1)  
    )  
  
    # Backward pass  
    loss.backward()  
    optimizer.step()
```

What it measures:

- How well model predicts next word given previous words
- Uses ground truth words as input (teacher forcing)
- Optimized during training

Validation Loss

```
python
```

```
model.eval()  
with torch.no_grad():  
    for mel, captions in val_loader:  
        logits = model(mel, captions[:, :-1])  
        targets = captions[:, 1:]  
  
        loss = criterion(  
            logits.reshape(-1, vocab_size),  
            targets.reshape(-1)  
        )  
  
        val_loss += loss.item()
```

What it measures:

- Same as training loss
- But on unseen data
- No gradient updates
- Tells us if model generalizes

Why Both Matter

Training Loss ↓ but Validation Loss ↑

```
Epoch 1: Train=2.5, Val=2.6 ✓
Epoch 5: Train=1.2, Val=1.4 ✓
Epoch 10: Train=0.5, Val=1.8 ✗ OVERFITTING!
    └ Model memorizing training data
```

Both Decrease:

```
Epoch 1: Train=2.5, Val=2.6 ✓
Epoch 5: Train=1.2, Val=1.3 ✓
Epoch 10: Train=0.5, Val=0.6 ✓ Good generalization!
```

Why Cross-Entropy Loss?

Perfect Prediction

```
python
Target: "engine" (id=145)
Model output: [0.0, 0.0, 1.0, 0.0, ...] (100% confident)
    index=145 ↑

Loss = -log(1.0) = 0.0 ← PERFECT!
```

Terrible Prediction

```
python
Target: "engine" (id=145)
Model output: [0.0, 0.0, 0.001, 0.0, ...] (0.1% confident)
    index=145 ↑

Loss = -log(0.001) = 6.9 ← TERRIBLE!
```

Uncertainty Penalty

```
python
```

```
# Model should be confident!
```

```
Confident: [0.0, 0.9, 0.05, 0.05] → Loss = 0.11
```

```
Uncertain: [0.25, 0.25, 0.25, 0.25] → Loss = 1.39
```

Even **if** correct word has highest probability,
being uncertain (flat distribution) **is** penalized!

Loss During Generation (Inference)

Key Difference

Training: Use ground truth words

```
Input: <sos> → an → engine → is
```

```
Target: an → engine → is → running
```

Inference: Use model's own predictions

```
Input: <sos> → [generate "an"] → [generate "engine"] → [generate "is"]
```

No loss calculated during inference!

Why? We don't have ground truth captions for new audio.

Other Evaluation Metrics

Loss is useful during training, but we need better metrics for final evaluation:

1. Perplexity

```
python
```

```
perplexity = exp(average_loss)
```

```
# Lower is better
```

```
Loss = 0.5 → Perplexity = 1.65 (very good)
```

```
Loss = 2.0 → Perplexity = 7.39 (okay)
```

```
Loss = 4.0 → Perplexity = 54.6 (bad)
```

Interpretation: On average, model is confused between ~perplexity choices

2. BLEU Score

```
python
```

```
from nltk.translate.bleu_score import sentence_bleu

reference = "an engine is running loudly"
generated = "a motor is running"

bleu = sentence_bleu([reference.split()], generated.split())
# BLEU = 0.58 (0-1 scale, higher is better)
```

Measures: N-gram overlap with reference captions

3. METEOR, ROUGE, CIDEr

Different ways to compare generated vs. reference captions

4. Human Evaluation

Ultimate test: Do humans think captions are good?

Visualizing Loss

What different loss values mean:

```
python
```

```
Loss = 0.0 - 0.5: Excellent (model very confident)
Loss = 0.5 - 1.0: Good (model mostly right)
Loss = 1.0 - 2.0: Okay (model learning)
Loss = 2.0 - 4.0: Poor (model confused)
Loss > 4.0:       Terrible (model random)
```

During Training

```
Epoch 1: Loss=4.2 (basically random guessing)
Epoch 5: Loss=2.1 (learning common patterns)
Epoch 10: Loss=1.2 (decent predictions)
Epoch 20: Loss=0.6 (good captions)
Epoch 30: Loss=0.4 (very good captions)
```

Code Example: Manual Loss Calculation

```
python
```

```
import torch
import torch.nn.functional as F

# Simplified example
vocab_size = 100
seq_len = 5
batch_size = 2

# Model output (logits)
logits = torch.randn(batch_size, seq_len, vocab_size)
# Target caption (word IDs)
targets = torch.randint(0, vocab_size, (batch_size, seq_len))

print(f"Logits shape: {logits.shape}")
print(f"Targets shape: {targets.shape}")

# Method 1: Using PyTorch's CrossEntropyLoss
criterion = nn.CrossEntropyLoss()
loss = criterion(
    logits.view(-1, vocab_size), # [batch*seq_len, vocab_size]
    targets.view(-1)           # [batch*seq_len]
)
print(f"Loss: {loss.item():.4f}")

# Method 2: Manual calculation
# Convert logits to probabilities
probs = F.softmax(logits, dim=-1) # [batch, seq_len, vocab_size]

# Get probability of correct word at each position
correct_probs = probs.gather(
    dim=2,
    index=targets.unsqueeze(-1)
).squeeze(-1) # [batch, seq_len]

# Cross-entropy = -log(probability of correct word)
manual_loss = -torch.log(correct_probs).mean()
print(f"Manual loss: {manual_loss.item():.4f}")

# They should match!
assert torch.allclose(loss, manual_loss, atol=1e-5)
```

Common Questions

Q1: Why not use MSE (regression) loss?

Answer: We're predicting discrete categories (words), not continuous values.

```
python
```

```
# This makes no sense:  
predicted_word_id = 145.7 # "engine" is 145, but 145.7??  
actual_word_id = 145  
  
mse_loss = (145.7 - 145)^2 = 0.49 # What does this mean?
```

Q2: Why ignore padding tokens?

```
python
```

```
criterion = nn.CrossEntropyLoss(ignore_index=vocab['<pad>'])
```

Answer: Padding is artificial, we don't want model learning to predict padding!

```
python
```

```
# Caption: "a sound" (short, needs padding)  
Actual: ["a", "sound", "<eos>", "<pad>", "<pad>", "<pad>"]  
       ↑ Don't penalize predicting these!
```

Q3: What about label smoothing?

```
python
```

```
criterion = nn.CrossEntropyLoss(label_smoothing=0.1)
```

Instead of:

```
python
```

```
Target: "engine"  
Hard target: [0, 0, 1.0, 0, 0, ...] # 100% confident
```

Use:

```
python
```

```
Soft target: [0.02, 0.02, 0.92, 0.02, 0.02, ...] # 92% confident  
          ↑ Small probability everywhere
```

Why? Prevents overconfidence, better generalization.

Q4: Why is validation loss sometimes lower than training?

Reasons:

1. **Dropout disabled during validation** (model more powerful)
2. **Batch size effects** (validation might have easier batches)
3. **Label smoothing** only applied during training

Not necessarily a problem if difference is small!

Summary Table

Aspect	Details
Problem Type	Classification (4,374-way at each position)
Loss Function	Cross-Entropy Loss
Formula	$-\log P(\text{correct_word})$
Range	0 (perfect) to ∞ (terrible)
Good Value	< 1.0 for audio captioning
Training	Uses ground truth words (teacher forcing)
Validation	Same as training, but on unseen data
Inference	No loss (no ground truth available)

The Bottom Line

Caption generation = Sequence of classification problems

Each word position:

- Model outputs probability distribution over 4,374 words
- Cross-entropy loss measures how confident model is in correct word
- Lower loss = better predictions
- Training minimizes this loss
- Validation checks generalization

Loss of 0.4-0.6 is typically good for audio captioning!

Your model's loss tells you:

- How well it predicts words
- Whether it's overfitting (train vs. val gap)
- When to stop training (early stopping)

But remember: **Low loss ≠ good captions necessarily**

Always also evaluate with:

- Generated sample captions (qualitative)
- BLEU/METEOR scores (quantitative)
- Human evaluation (ultimate test)