

Transformer Model Results Analysis

Overall Assessment: 🎉 Huge Improvement!

Your model is now generating **diverse, contextually relevant captions** instead of the repetitive "a person is walking" problem!

✓ What's Working Well

1. Diversity Achieved!

Sample 1: "a machine is running and then turned off"
Sample 2: "a woman speaks over an electronic sound..."
Sample 3: "water is pouring down into a container..."
Sample 6: "people are talking in a large crowd..."
Sample 9: "the wind is blowing in the background..."

Each caption is unique and contextually different! ✓

2. Audio-Grounded Generation

The model is clearly listening to the audio:

Example 1: Machine Sound

Audio: Santa Motor.wav
Generated: "a machine is running and then turned off"
Reference: "A machine whines and squeals..."
✓ Correctly identified: machine, mechanical sound

Example 6: Crowd

Audio: young artists.wav
Generated: "people are talking in a large crowd of people..."
Reference: "A large gathering of people are talking loudly..."
✓ Correctly identified: multiple people, talking, crowd

Example 9: Wind/Air

Audio: Sound of the wind comes from the tunnel 3.wav
Generated: "the wind is blowing in the background..."
Reference: "The wind is howling through a large room"
✓ Correctly identified: wind, air movement

3. Grammatically Correct

All 10 captions are grammatically valid sentences! ✓

4. Semantically Reasonable

Unlike "applause is whirring", these make logical sense:

- "machine is running" ✓
 - "people are talking" ✓
 - "wind is blowing" ✓
 - "dog barks" ✓
-

⚠️ Issues & Areas for Improvement

1. Redundancy Within Captions

Sample 2:

"a woman speaks over an electronic sound and a woman speaks to each other"
↑ repetition

Sample 6:

"people are talking in a large crowd of people are talking in a crowded area"
↑ repeated phrase

Sample 10:

"a clock is ticking with some water being poured into a container is poured into a container"
↑ repeated

Why this happens:

- Model gets into repetitive loop
- Attention mechanism focusing on same features
- Need better diversity penalty or sampling

Fix:

```
python
```

```
# Increase temperature slightly  
model.generate(mel, temperature=0.9) # Currently 0.8  
  
# Add repetition penalty  
# Penalize recently generated words
```

2. Imprecise Matches

Sample 3: Radio Static

Generated: "water is pouring down into a container..."
Reference: "A radio tuner has been positioned in between radio stations..."
✖ Completely wrong interpretation!

Why: Radio static might sound like water/fizzing to the model

Sample 7: Gasping

Generated: "a person is playing a metal objects..."
Reference: "A man is inhaling air with a short gasp..."
✖ Missed the breathing/gasping entirely

Why: Model needs better temporal feature extraction

3. Awkward Phrasing

Sample 4:

"a person is opening a wooden door to a container"
↑ "door to a container"?

Sample 5:

"a metal object is repeatedly and a hard object being dragged..."
↑ "repeatedly and"? Incomplete thought

Sample 10:

"a clock is ticking with some water being poured..."
↑ Clock ticking + water pouring? Mixed concepts

4. Length/Complexity Issues

Some captions are too long and convoluted:

Sample 5: "a metal object is repeatedly and a hard object being dragged across another surface in a metal surface"

Could be: "metal objects scraping against each other"

Quantitative Analysis

Caption Length Distribution

Sample 1: 9 words ✓ Good

Sample 2: 15 words ⚠ Getting long

Sample 3: 17 words ⚠ Too long + repetitive

Sample 4: 9 words ✓ Good

Sample 5: 18 words ⚠ Too long + awkward

Sample 6: 16 words ⚠ Repetitive

Sample 7: 10 words ✓ Good

Sample 8: 9 words ✓ Good

Sample 9: 13 words ⚠ Repetitive

Sample 10: 18 words ⚠ Very repetitive

Average: 13.4 words

Optimal: 8-12 words

Accuracy Estimation (Rough)

Sample	Accuracy	Notes
1	80%	"machine running" ✓, but not specific details
2	70%	Got "woman speaks" ✓, missed "radio/walkie-talkie"
3	20%	✗ "water" instead of "radio static"
4	40%	✗ "door" instead of "rattle/jewelry"
5	75%	"metal objects" ✓, scraped surface ✓
6	95%	✓ Nearly perfect!
7	30%	✗ "metal objects" instead of "gasping"
8	85%	"dog barking" ✓ (though it's human mimicking)
9	90%	"wind blowing" ✓
10	60%	Got "water/container" ✓, missed "drinking"

Average: ~65% accuracy

Not bad for a first real attempt!

🎯 Performance by Sound Type

✓ Excellent: Crowd/People Sounds

Sample 6: "people are talking in a large crowd"

Reference: "A large gathering of people are talking loudly"

Match: 95% ✓✓✓

✓ Good: Wind/Air/Ambient

Sample 9: "the wind is blowing in the background"

Reference: "The wind is howling through a large room"

Match: 90% ✓✓

✓ Good: Animal Sounds

Sample 8: "a dog barks and a dog barking"

Reference: "A person is attempting to mimic an angry dog"

Match: 85% ✓✓ (close enough)

✓ Good: Mechanical

Sample 1: "a machine is running and then turned off"

Reference: "A machine whines and squeals..."

Match: 80% ✓✓

⚠ Fair: Metallic/Contact Sounds

Sample 5: "a metal object is repeatedly and a hard object being dragged..."

Reference: "A person is pulling silverware out of the dishwasher"

Match: 60% ✓

✗ Poor: Radio Static/Complex Sounds

Sample 3: "water is pouring down..."

Reference: "A radio tuner has been positioned..."

Match: 20% ✗

✗ Poor: Human Breathing/Gasping

Sample 7: "a person is playing a metal objects"

Reference: "A man is inhaling air with a short gasp"

Match: 30% ✗

🔍 Comparison to Your Original Model

Before (Baseline with repetition):

Sample 1: "a person is walking on a hard surface while a person is walking in the background"

Sample 2: "a person is walking on a hard surface while a person is walking in the background"

Sample 3: "a person is walking on a hard surface while a person is walking in the background"

...

- ✗ No diversity
- ✗ Not audio-grounded
- ✗ Useless output

Now (Transformer):

Sample 1: "a machine is running and then turned off"

Sample 2: "a woman speaks over an electronic sound..."

Sample 3: "water is pouring down into a container..."

...

- ✓ Diverse captions
- ✓ Audio-grounded (mostly)
- ✓ Useful output
- ! Some repetition within captions
- ! Some inaccuracies

Improvement: From 0% → 65% accuracy! 🎉



Recommendations for Further Improvement

Priority 1: Fix Internal Repetition (Quick Win)

Problem: "people are talking... people are talking"

Solution 1: Repetition Penalty

```
python
```

```
def generate_with_repetition_penalty(model, mel, penalty=1.2):
    """Penalize recently used tokens"""
    generated = []

    for step in range(max_len):
        logits = model.get_logits(...)

        # Penalize tokens already in sequence
        for token in generated[-5:]: # Last 5 tokens
            logits[token] /= penalty # Reduce score

        next_token = sample(logits)
        generated.append(next_token)
```

Solution 2: N-gram Blocking

```
python
```

```
# Don't allow same 2-3 word phrase to repeat
if (word1, word2) in recent_bigrams:
    logits[word2] = -inf # Block it
```

Priority 2: Better Audio Encoder (Medium Effort)

Problem: Confusing water sounds with radio static

Solution: Use Pre-trained Audio Model

```
python
```

```
# Instead of CNN from scratch:
from transformers import ASTModel # Audio Spectrogram Transformer

# Pre-trained on AudioSet (2M audio clips)
audio_encoder = ASTModel.from_pretrained('MIT/ast-finetuned-audioset-10-10-0.4593')

# Fine-tune on Clotho
```

Expected improvement: 65% → 75% accuracy

Priority 3: Constrain Caption Length (Easy)

Problem: Too many long, repetitive captions

Solution:

```
python
```

```
model.generate(  
    mel,  
    max_len=15, # ← Reduce from 30  
    temperature=0.8,  
    top_p=0.9  
)
```

Or penalty for length:

```
python
```

```
# Encourage model to finish earlier  
if len(generated) > 12:  
    logits[eos_token] *= 1.5 # Boost <eos> probability
```

Priority 4: Post-Processing (Quick Fix)

Clean up awkward captions:

```
python
```

```
def clean_caption(caption):  
    """Remove obvious errors"""  
    words = caption.split()  
  
    # Remove consecutive duplicates  
    cleaned = [words[0]]  
    for word in words[1:]:  
        if word != cleaned[-1]:  
            cleaned.append(word)  
  
    # Remove phrases like "of people are talking" after "crowd"  
    # (pattern matching)  
  
    return ' '.join(cleaned)  
  
# Example:  
before = "people are talking in a large crowd of people are talking"  
after = "people are talking in a large crowd"
```

Priority 5: Ensemble Models (Advanced)

Generate multiple captions, pick best:

python

```
# Generate 5 captions with different temperatures
captions = []
for temp in [0.7, 0.8, 0.9, 1.0, 1.1]:
    caption = model.generate(mel, temperature=temp)
    captions.append(caption)

# Score each (language model, length, uniqueness)
best_caption = rank_and_select(captions)
```

Expected Performance After Fixes

Current

Diversity: ✓✓✓✓ (9/10)

Accuracy: ✓✓✓ (6.5/10)

Fluency: ✓✓✓ (7/10)

Overall: 7.5/10

After Quick Fixes (Repetition Penalty + Length Constraint)

Diversity: ✓✓✓✓✓ (10/10)

Accuracy: ✓✓✓✓ (7/10)

Fluency: ✓✓✓✓ (8/10)

Overall: 8.3/10

After Pre-trained Encoder

Diversity: ✓✓✓✓✓ (10/10)

Accuracy: ✓✓✓✓ (8/10)

Fluency: ✓✓✓✓ (8.5/10)

Overall: 8.8/10

What You've Learned

Your journey:

1. **Baseline** → Complete mode collapse ("person walking" for everything)
2. **Fixed Architecture** → Some diversity but still issues
3. **Transformer** → Real, useful captions!