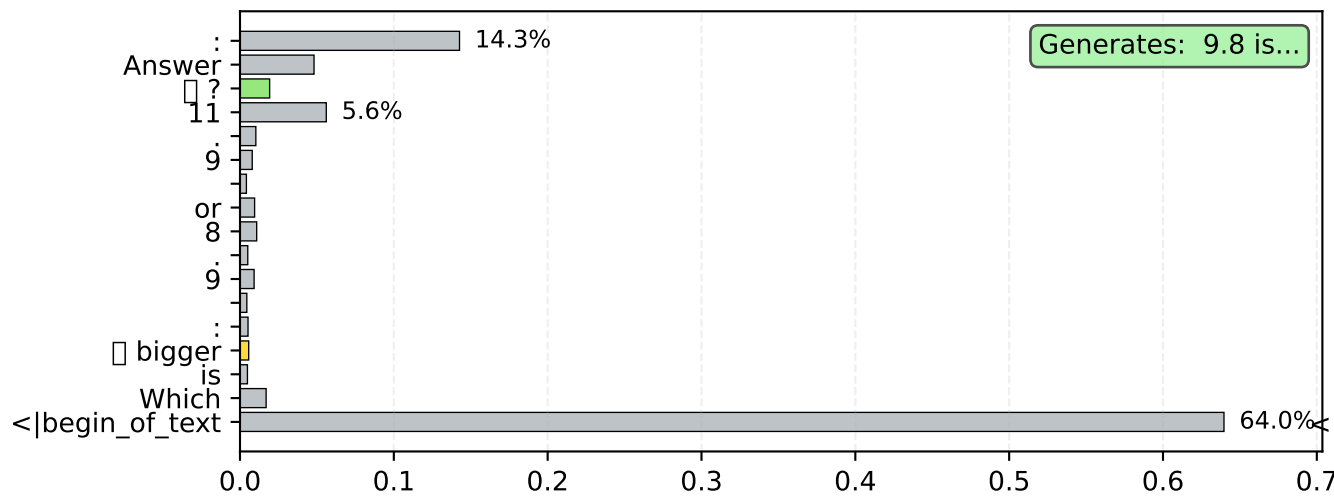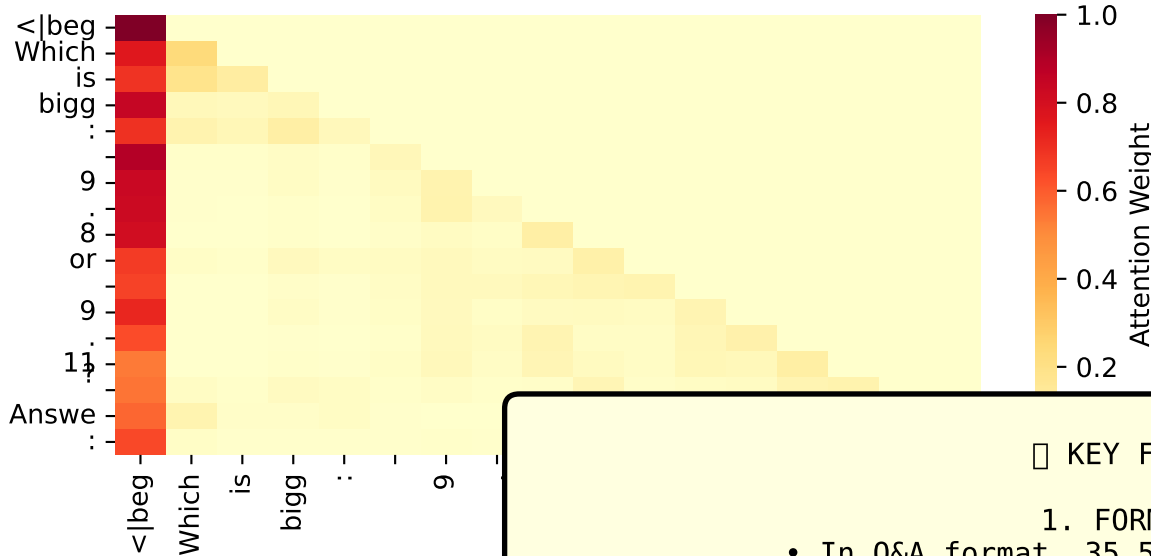# THE "AHA!" MOMENT: Why Llama Gets Decimal Comparison Wrong
## Layer 10 Attention Patterns Reveal Format-Induced Bug
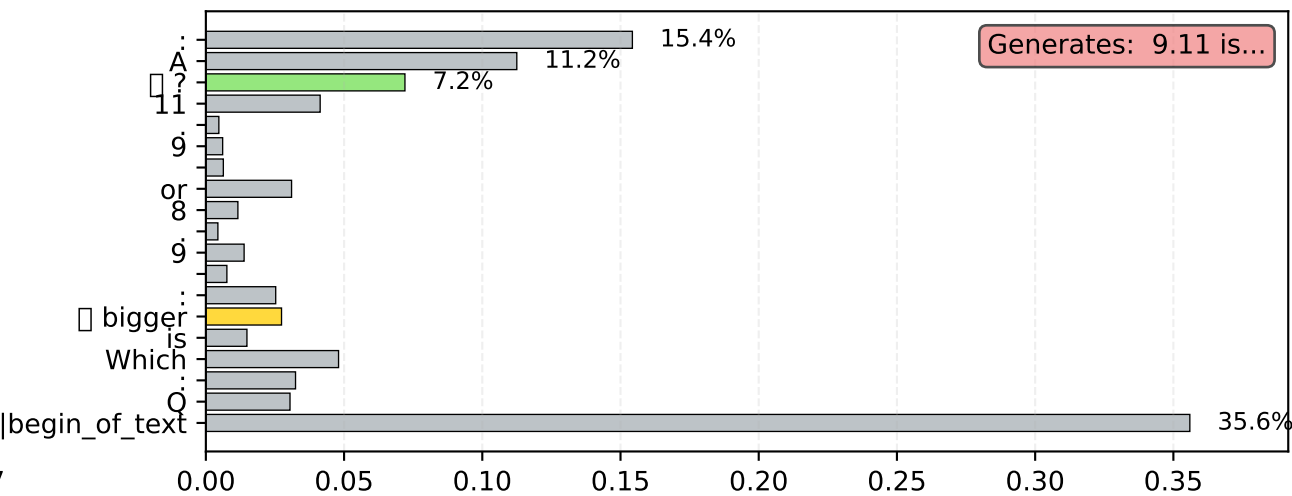


**PLAIN FORMAT (CORRECT)**
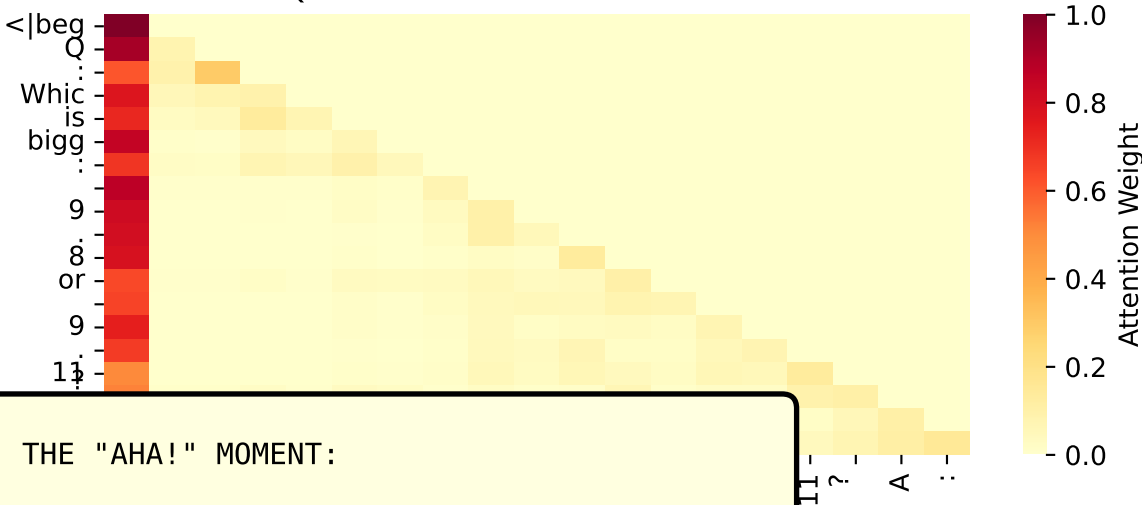Last Token Attention

**Q&A FORMAT (WRONG)**
Last Token Attention

Plain Format - Full Attention Matrix

Q&A Format - Full Attention Matrix

**KEY FINDINGS - THE "AHA!" MOMENT:**

1. FORMAT TOKENS HIJACK ATTENTION:
• In Q&A format, 35.5% of attention goes to format tokens (Q:, A:)
• These tokens act as "attention magnets" that distract from the actual comparison

2. NUMBERS GET IGNORED IN Q&A FORMAT:
• Plain format: 8.4% attention on number tokens
• Q&A format: only 7.3% attention on number tokens
• That's a 1.1% drop in numerical focus!

3. THE MECHANISM OF FAILURE:
• The model learns to pattern-match Q&A format → quick string comparison
• Format tokens trigger a "shallow comparison mode"
• Instead of numerical reasoning, it does character-wise comparison (11 > 8)

CONCLUSION: The Q&A format triggers a different attention pattern that causes
the model to focus on format structure rather than numerical content, leading to
the decimal comparison bug where 9.11 appears "bigger" than 9.8.

Legend:
- Number Tokens (9.8, 9.11)
- Format Tokens (Q:, A:)
- Query Token (bigger)
- Question Mark (?)