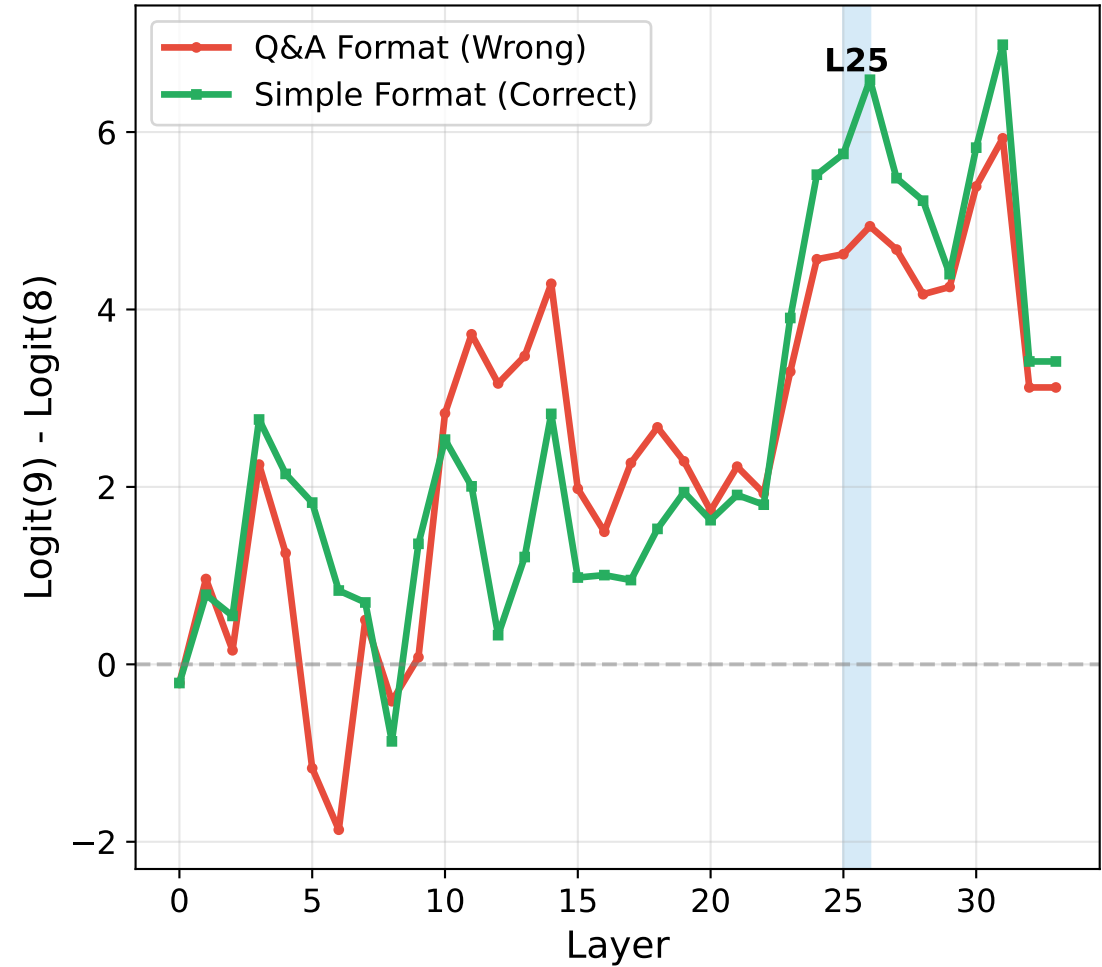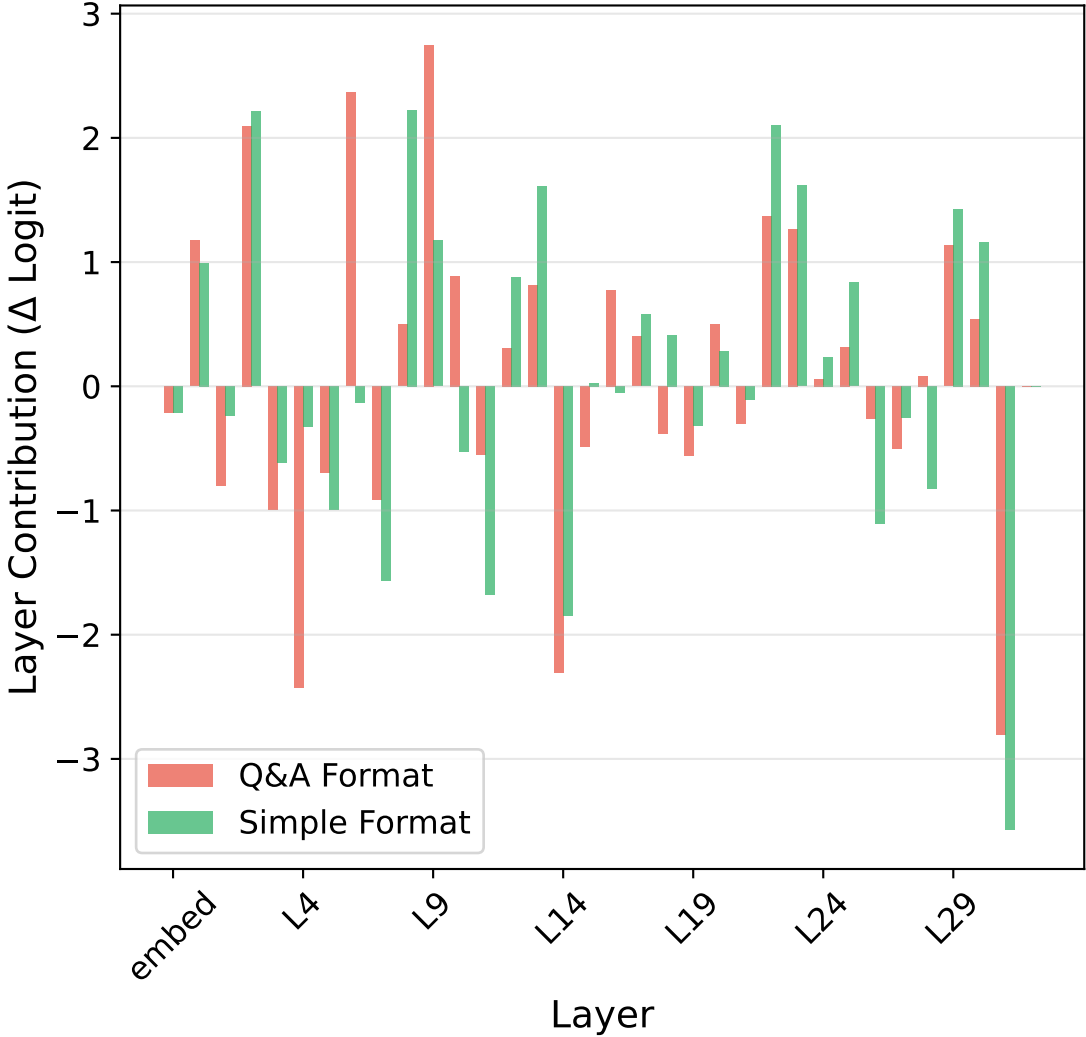Logit Attribution Analysis: Mechanistic Path of the Decimal Bug
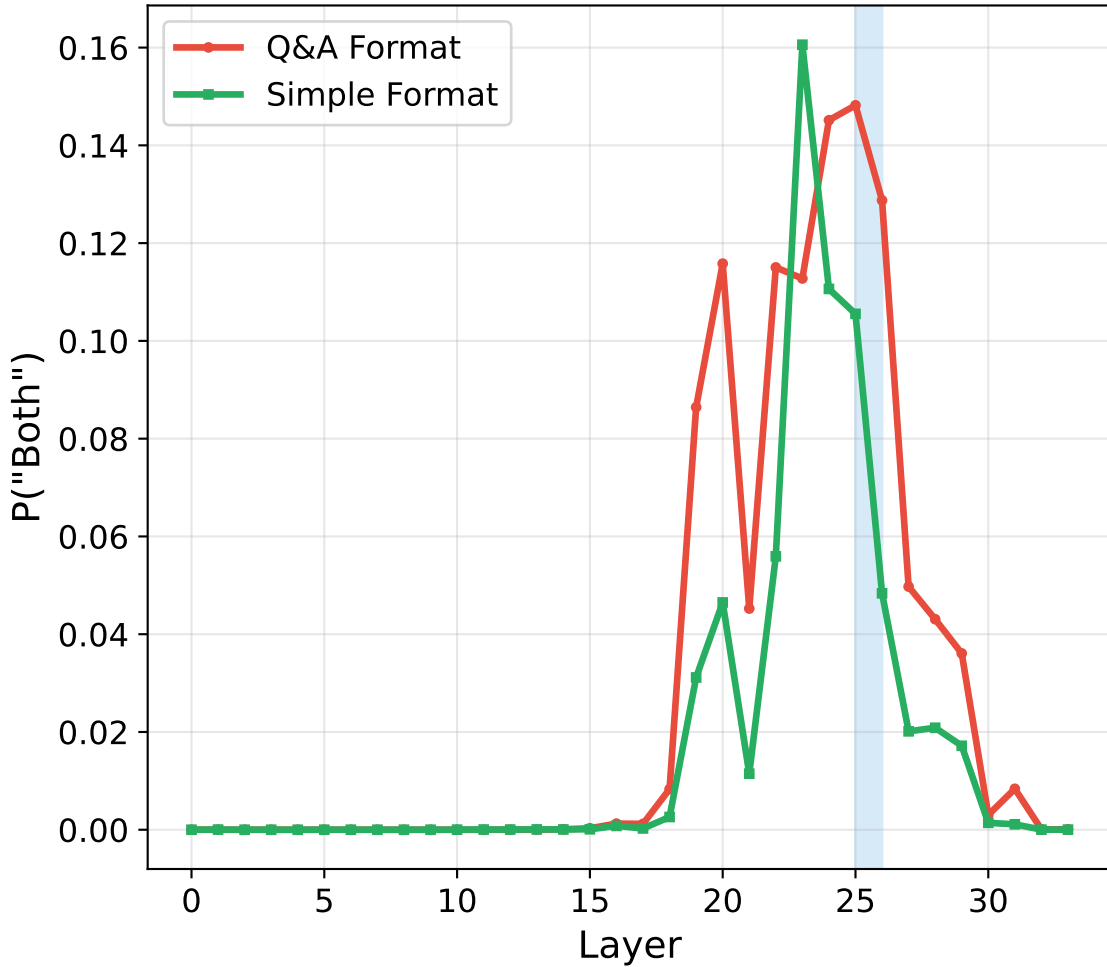
**Cumulative Logit Difference**
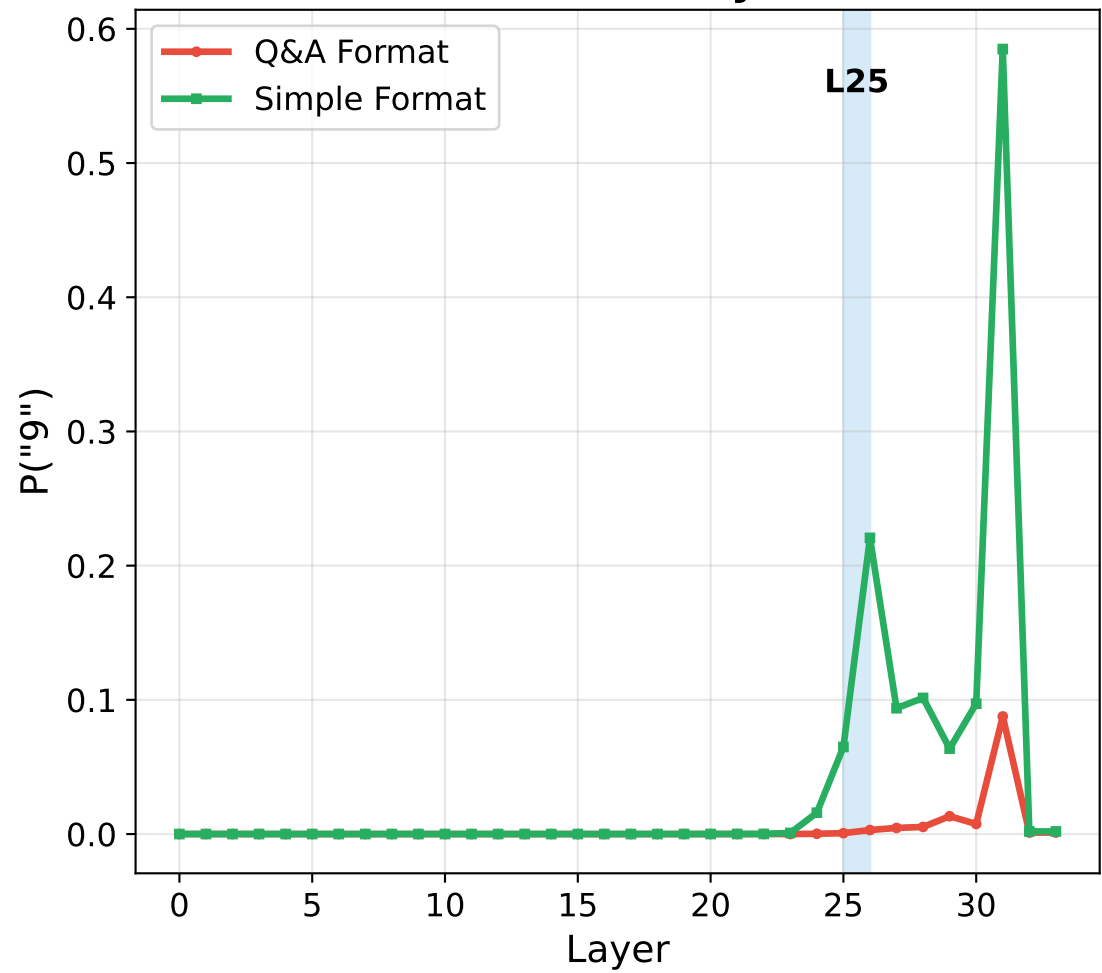
Q&A Format (Wrong)
Simple Format (Correct)

L25

**Per-Layer Contributions**

Q&A Format
Simple Format

**"Both" Token Probability**

Q&A Format
Simple Format

**Token "9" Probability (Critical)**

Q&A Format
Simple Format

L25

**Zoom: Layers 20-30 (Hedging Zone)**

Q&A Format
Simple Format

Layer 25
Divergence Point

KEY FINDINGS

Layer 25 Analysis:
• Q&A Format:
  - Δ Logit: 0.314
  - P("9"): 0.003
  - P("Both"): 0.129

• Simple Format:
  - Δ Logit: 0.836
  - P("9"): 0.221
  - P("Both"): 0.048

Critical Insight:
At Layer 25, Simple format
commits to "9" (22.1%)
while Q&A format hedges.