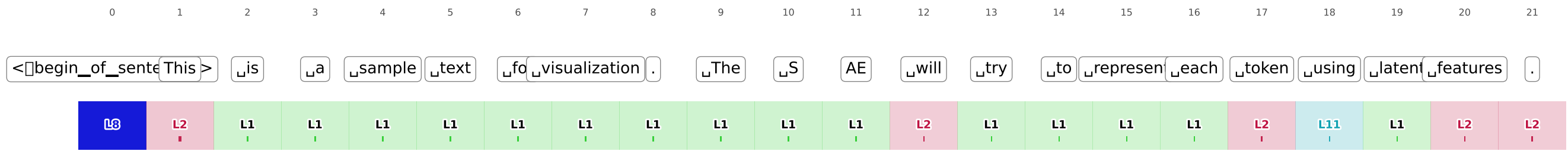# SAE Reasoning Trace - DEEPSEEK-R1-DISTILL-LLAMA-8B, Layer 6



Most Active Latent Features

- L1: Problem Ident...
- L8: Implication a...
- L11: Metacognitive...
- L2: Stepwise Nume...