Fig 1: NIAH task filled with repeative noise sentences
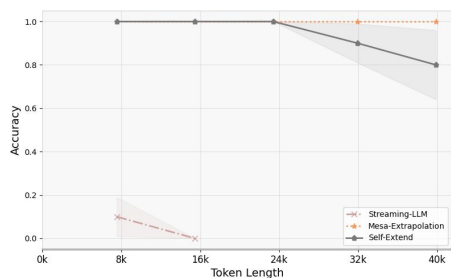

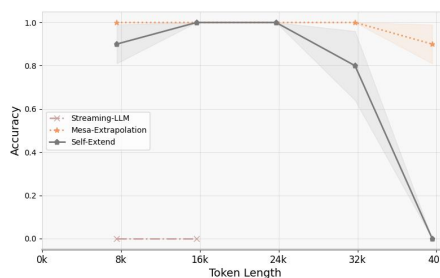Fig 2: NIAH task filled with Paul Graham essays


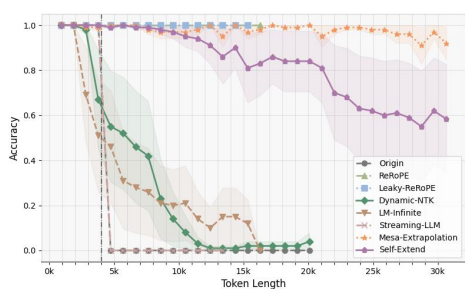Fig 3: refer from Figure 2 with added Self-Extend


Fig 4: NIAH tasks with single-multi keys using Phi3-mini-128k for Origin & ours
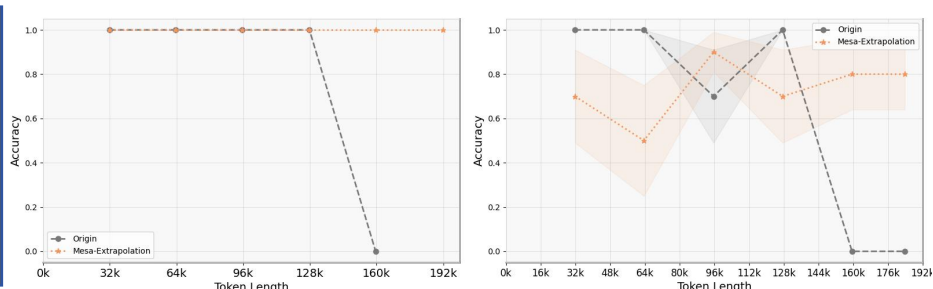

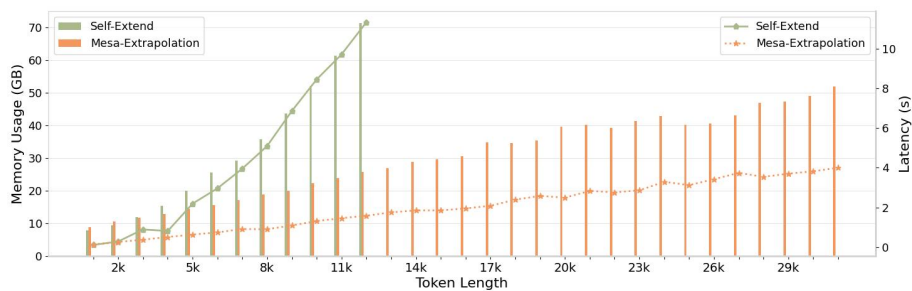Fig 5: Memory Usage & Latency on Open-LLaMA-3B model


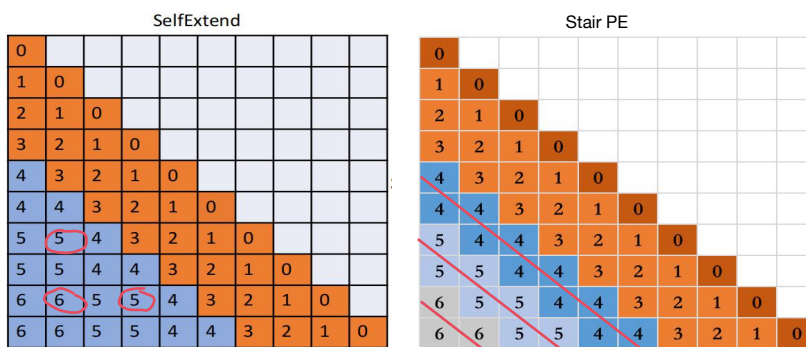Fig 6: attention matrix (before Softmax operation) compared between SelfExtend and Stair PE
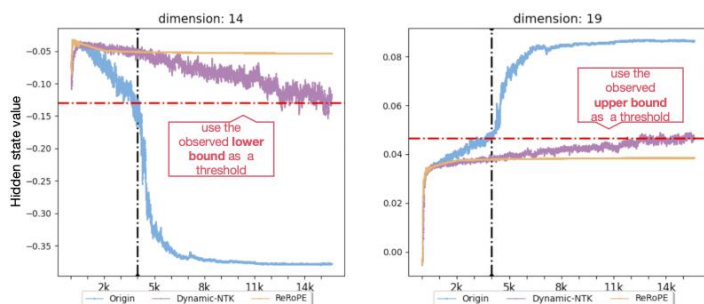

Fig 7: Dim 14 and 19 are observed with "good" probe


Fig 8: Dim 10 is observed with same passkey input.