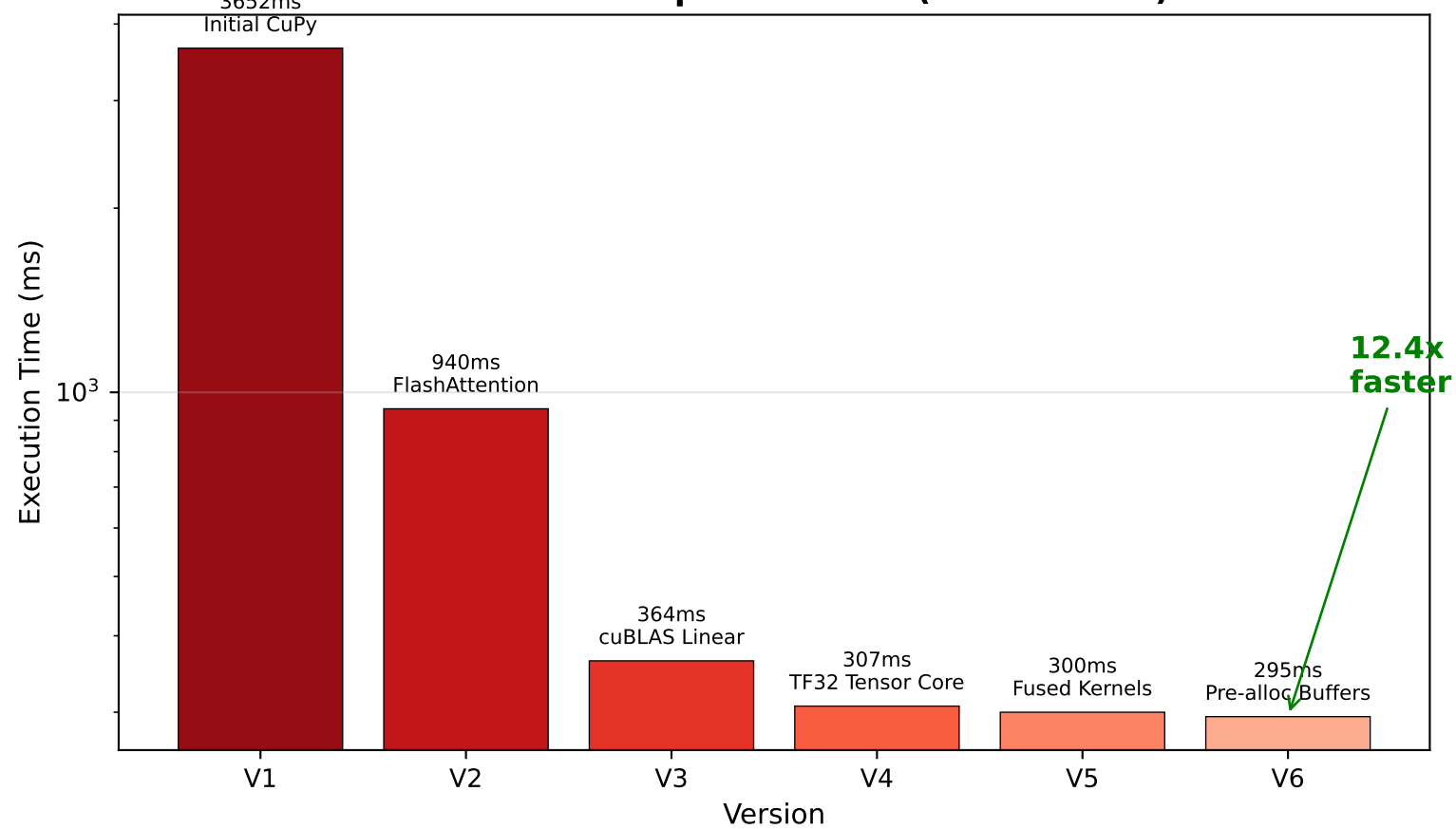
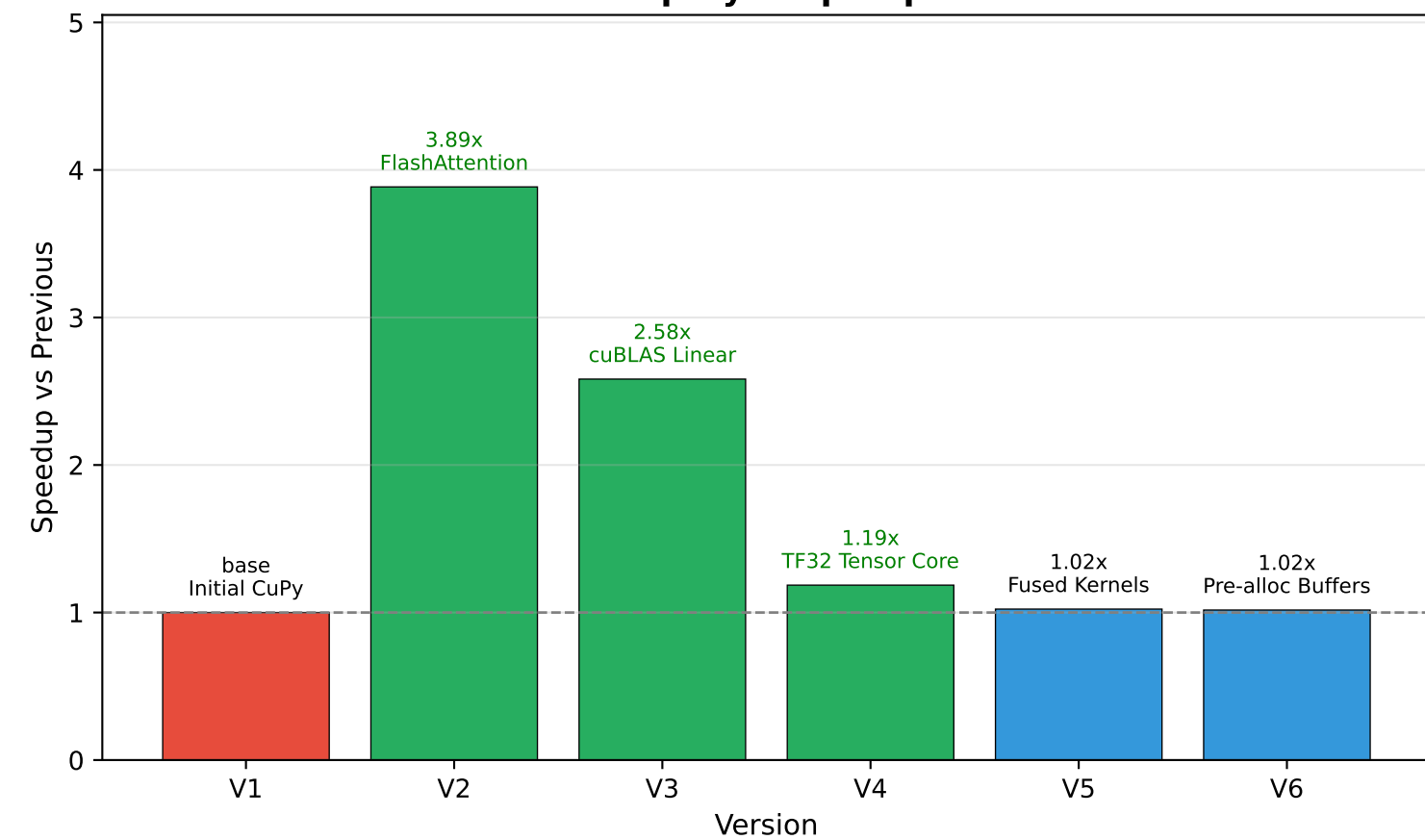


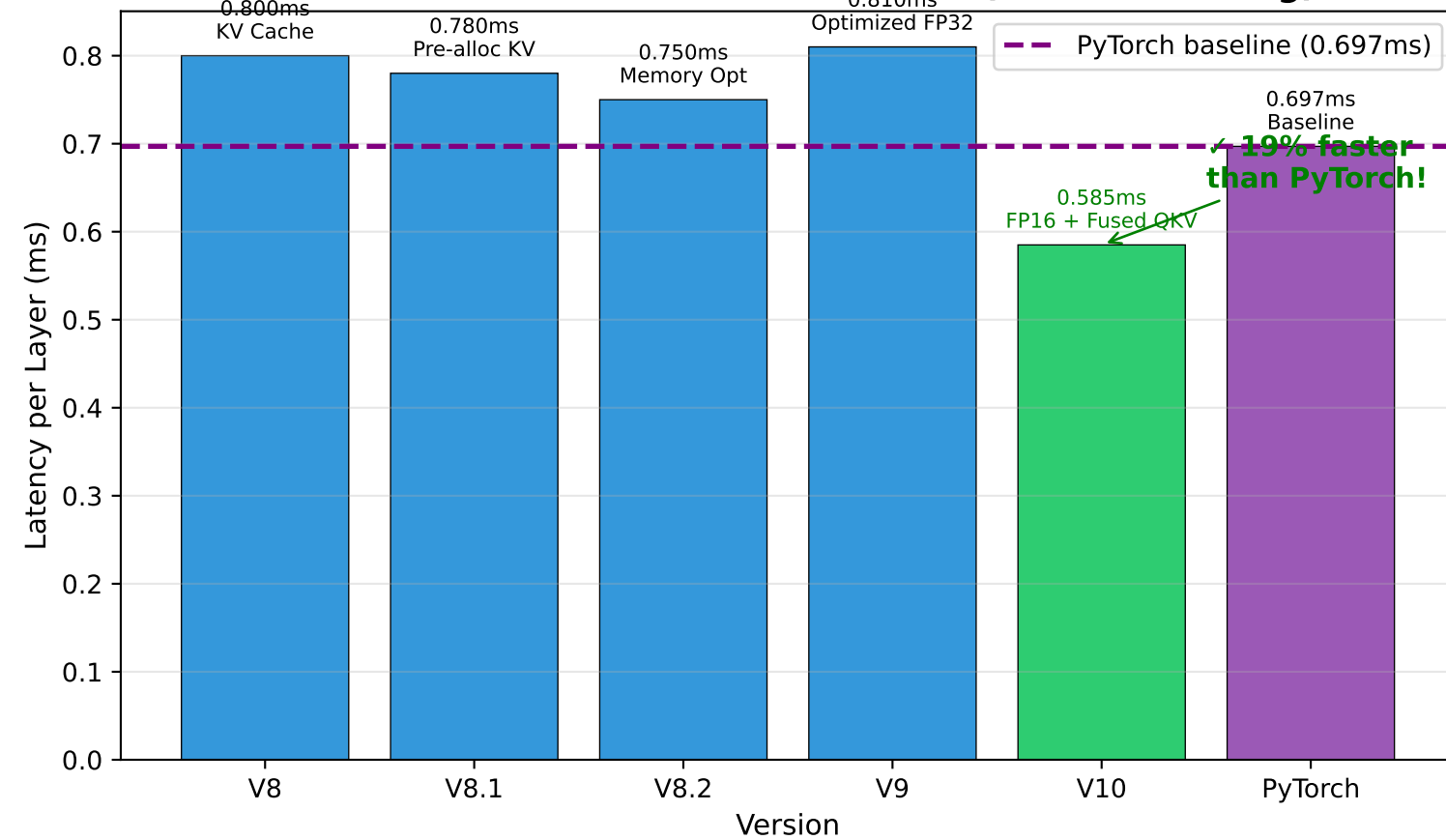
Phase 1: Core Optimizations (Prefill Mode)



Phase 1: Step-by-Step Improvement



Phase 2: Decode Mode with KV Cache (Real ASR Config)



Complete Optimization Journey: V1 → V10

