

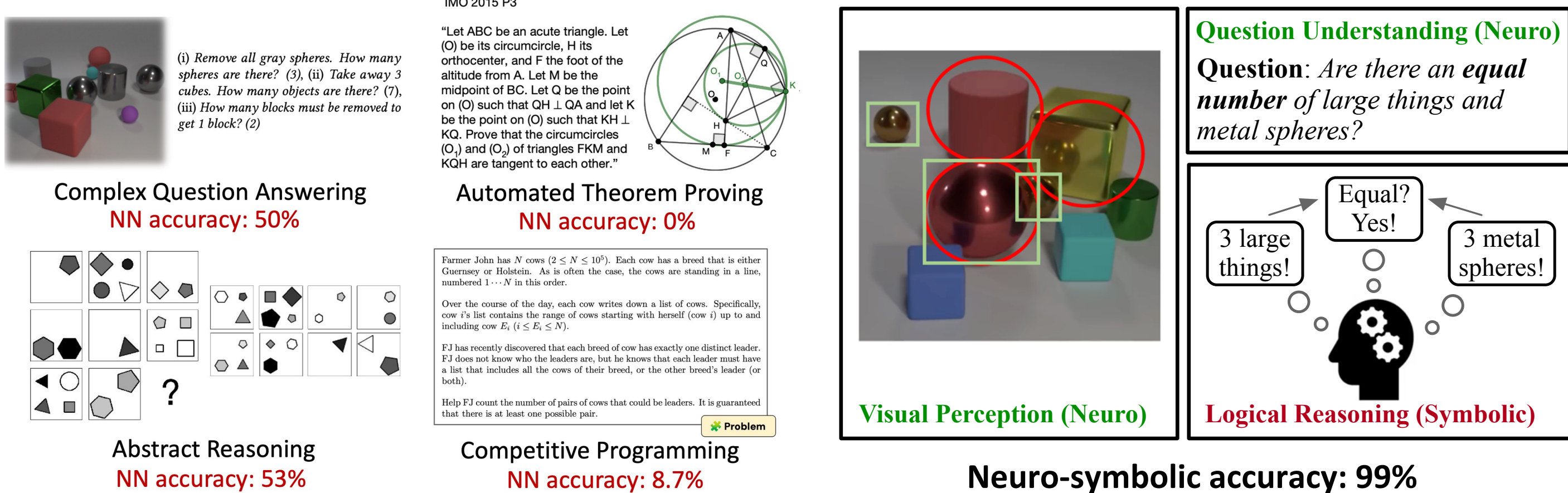
CogSys: Efficient and Scalable Neuro-Symbolic Cognition System via Algorithm-Hardware Co-Design

Zishen Wan^{1*}, Hanchen Yang^{1*}, Ritik Raj^{1*}, Che-Kai Liu¹, Anand Samajdar², Arijit Raychowdhury¹, Tushar Krishna¹

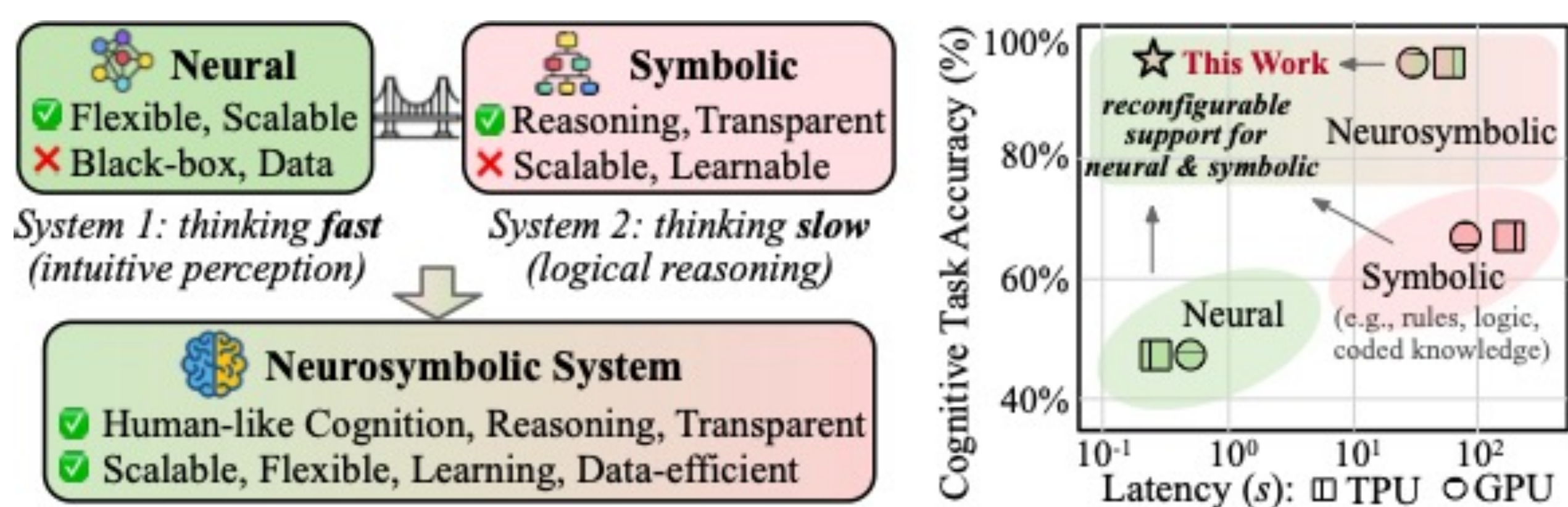
¹Georgia Tech, Atlanta, GA, USA ²IBM Research, Yorktown Heights, NY, USA

MOTIVATION: WHY NEURO-SYMBOLIC AI?

- Compositional system to enhance cognitive capability
- Applications: complex QA, abstract reasoning, math proving, etc

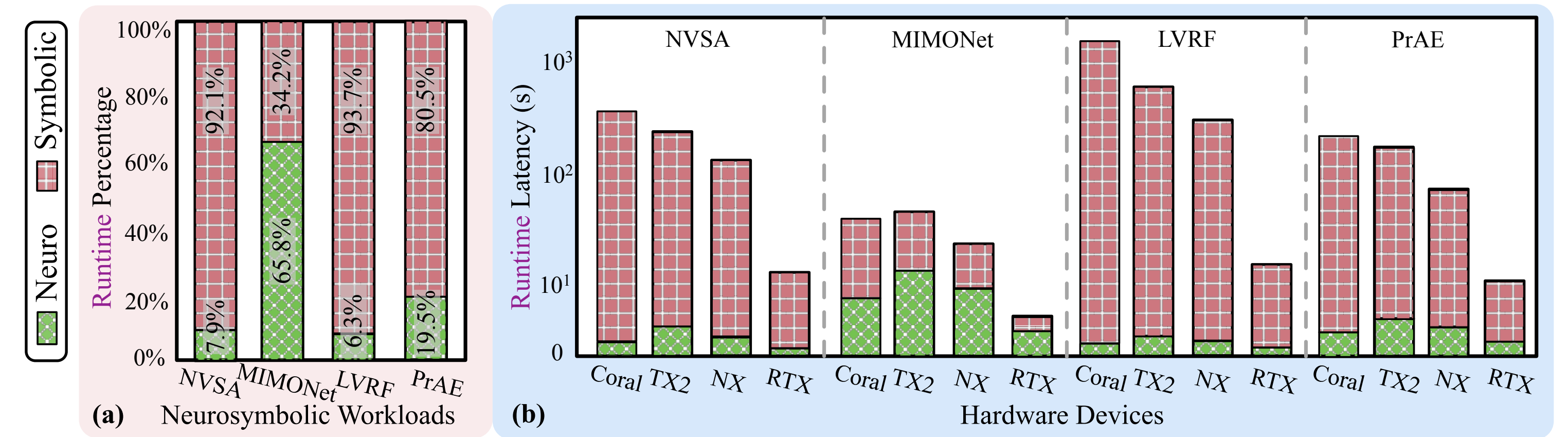


- Neuro-Symbolic AI bridges neural learning & symbolic reasoning

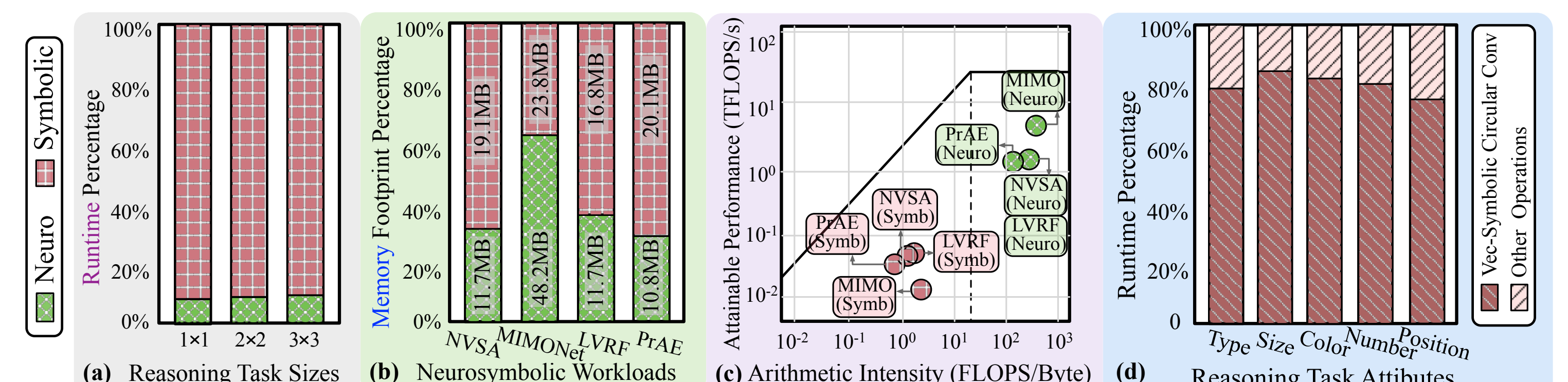


SYSTEM CHARACTERISTICS AND CHALLENGES

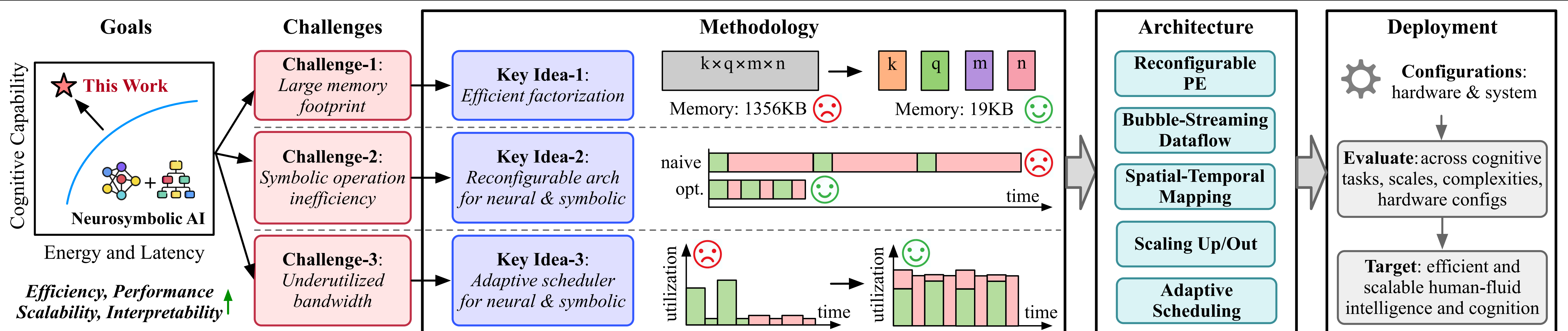
- System Challenges: Latency, memory, heterogeneity
- Latency: large end-to-end runtime; symbolic critical path bottleneck



- Memory: symbolic memory-bound; low ALU util, low cache hit rate
- Heterogeneity: neuro (GEMM, conv), symbolic (vector, circular conv)



METHODOLOGY: SOFTWARE-HARDWARE CO-DESIGN FOR REAL-TIME AND EFFICIENT NEURO-SYMBOLIC AI



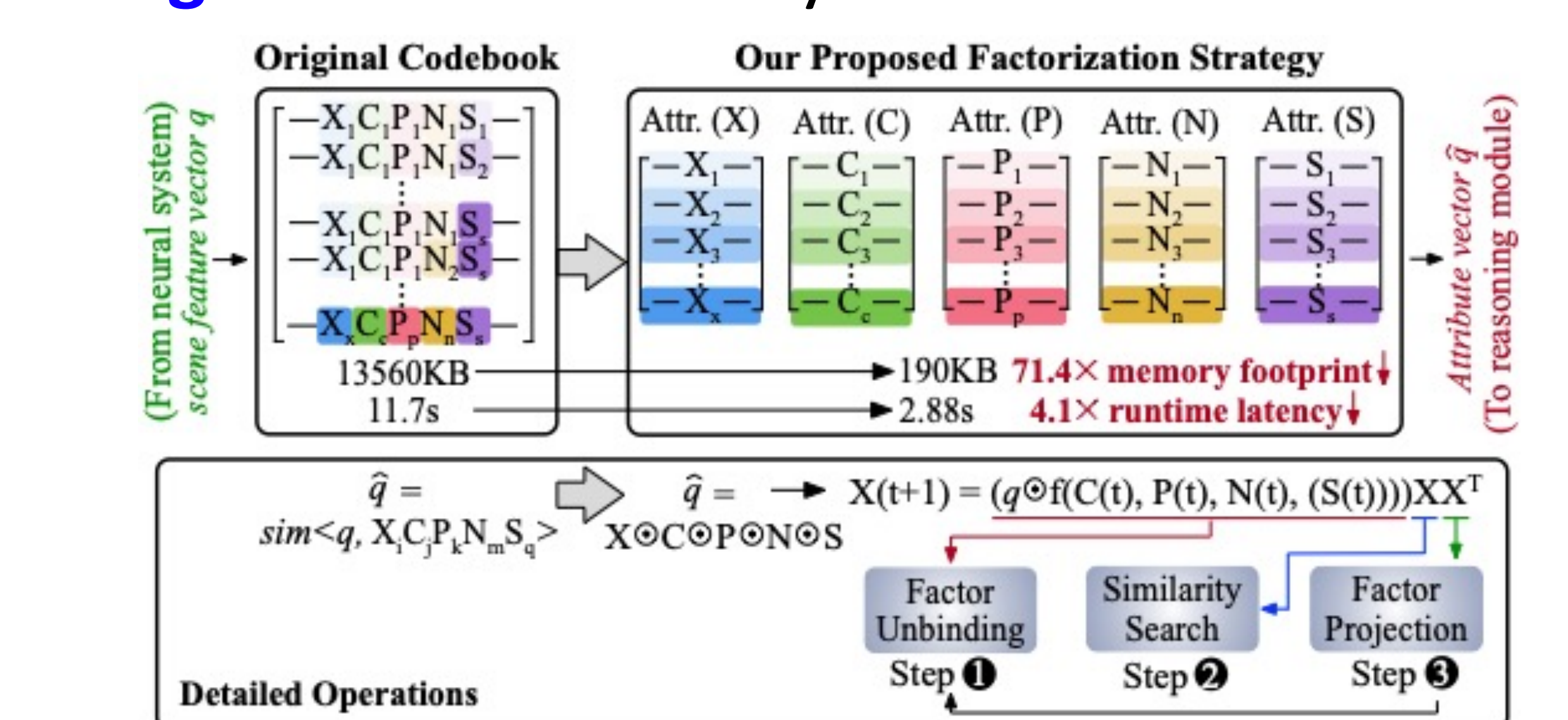
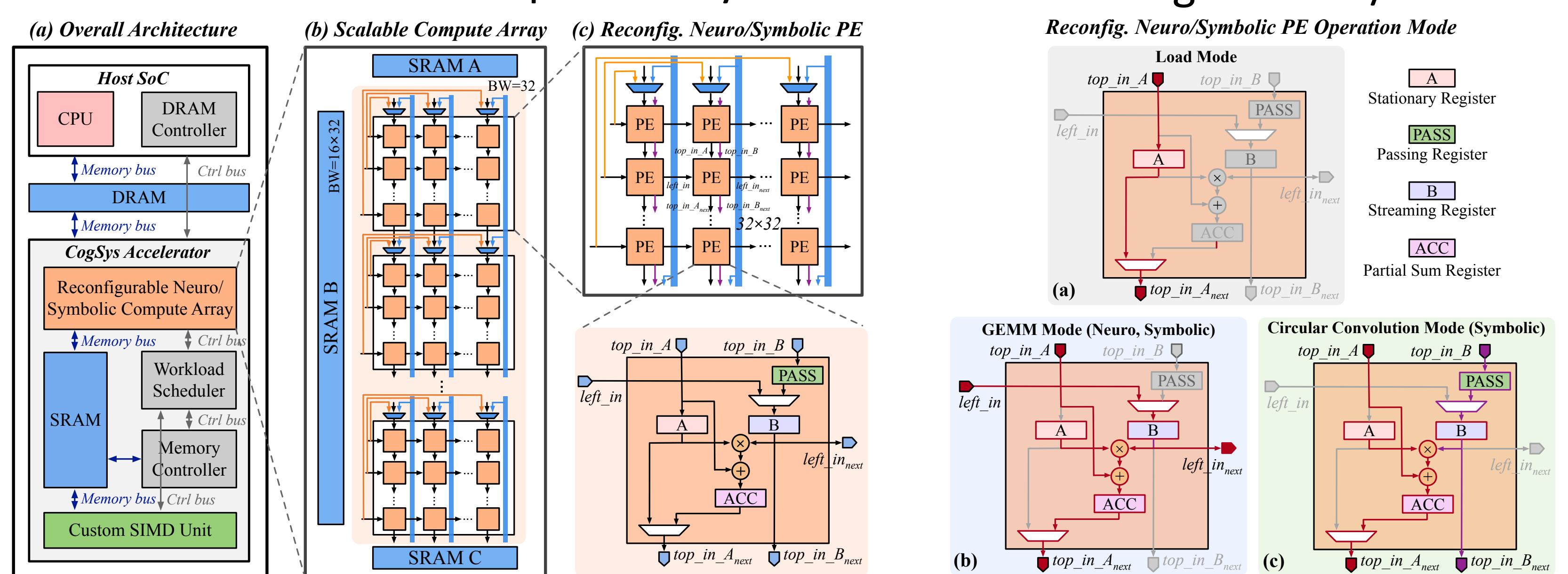
CROSS-LAYER OPTIMIZATION: ALGORITHM, HARDWARE ARCHITECTURE, SYSTEM SCHEDULING

- Hardware: Reconfigurable neuro-symbolic architecture

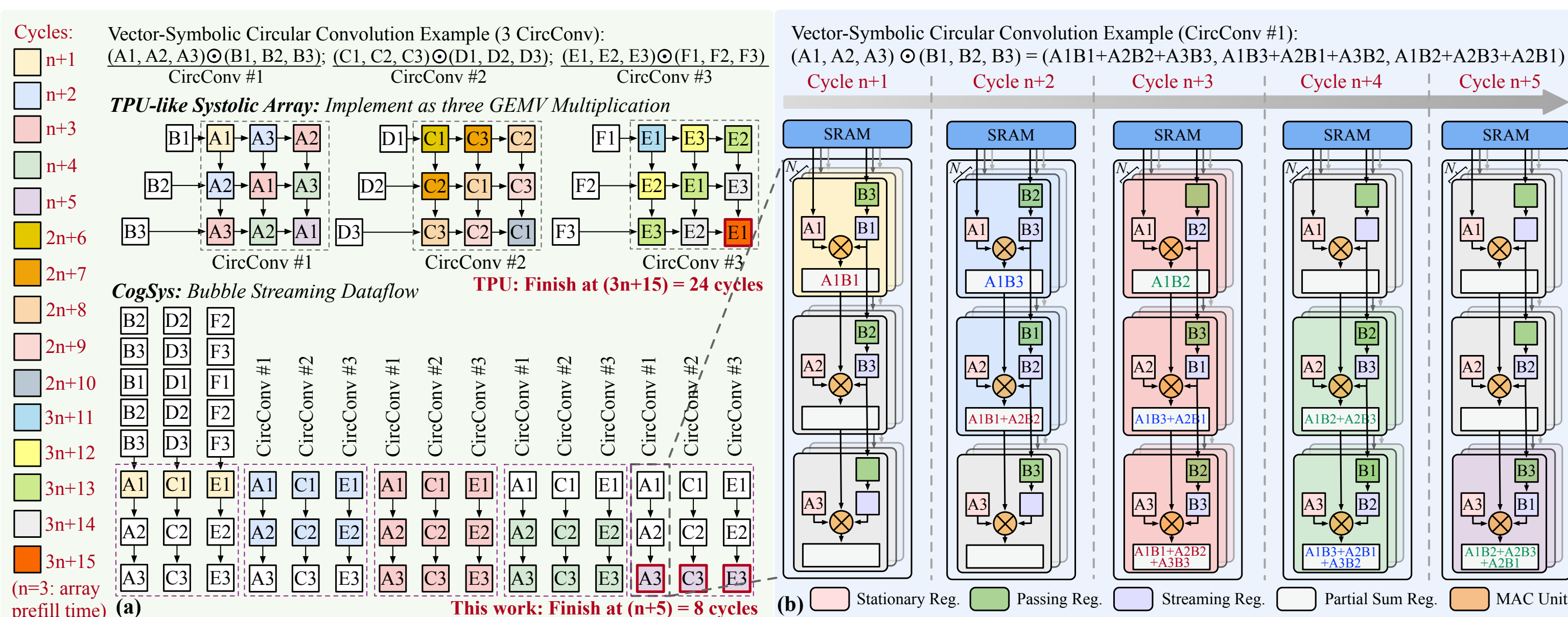
- Overview: scalable compute array

- PE: reconfigurable N/S modes

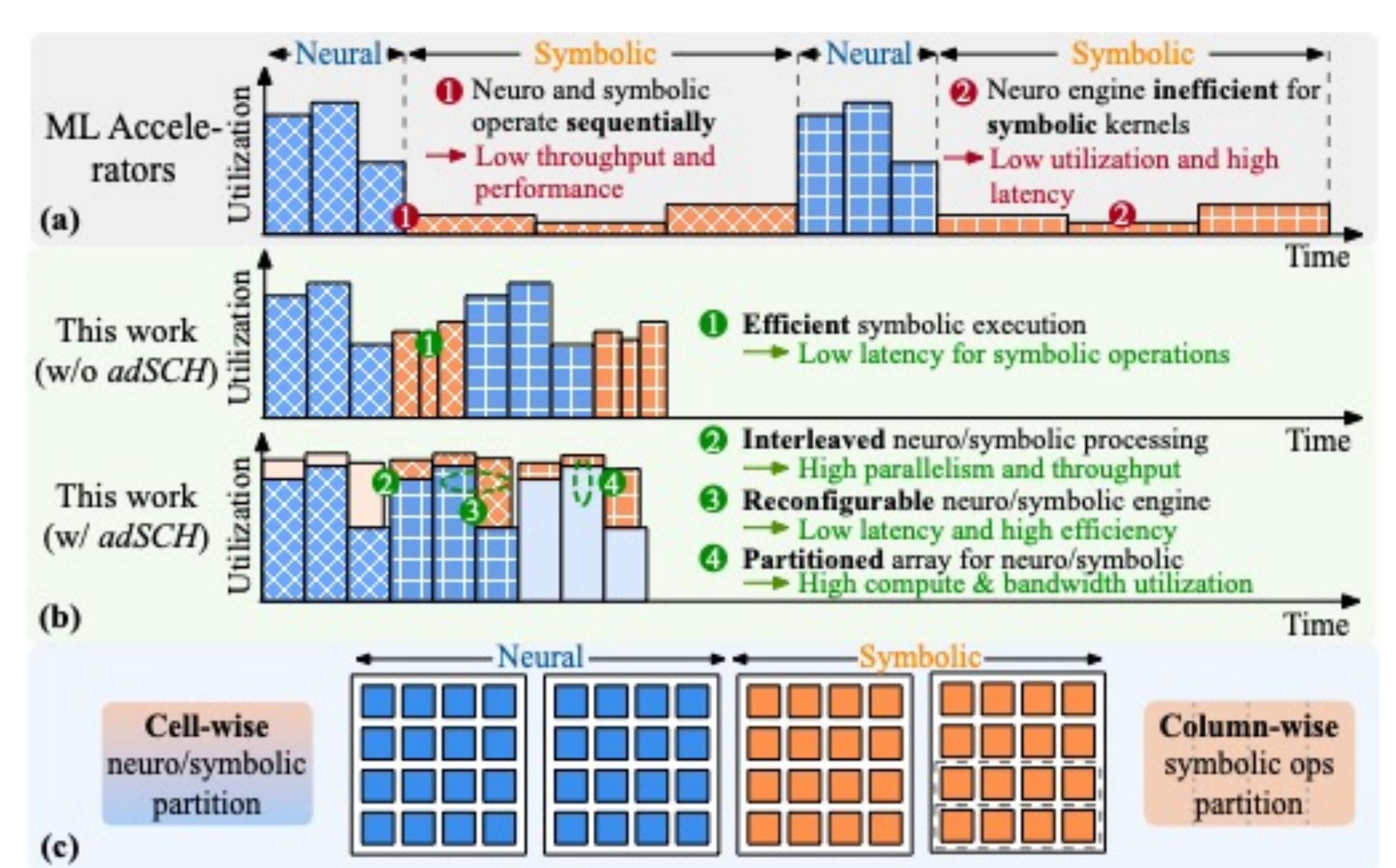
- Algorithm: Efficient symbolic factorization



- Dataflow: bubble streaming dataflow; adaptive spatial-temporal mapping



- System: Adaptive workload-aware scheduling



EVALUATION RESULTS

- Layout: 28nm node

- Latency & Energy: 90x speedup vs. edge GPU

- Compared with ML accelerators: 1.7-15.9x speedup

