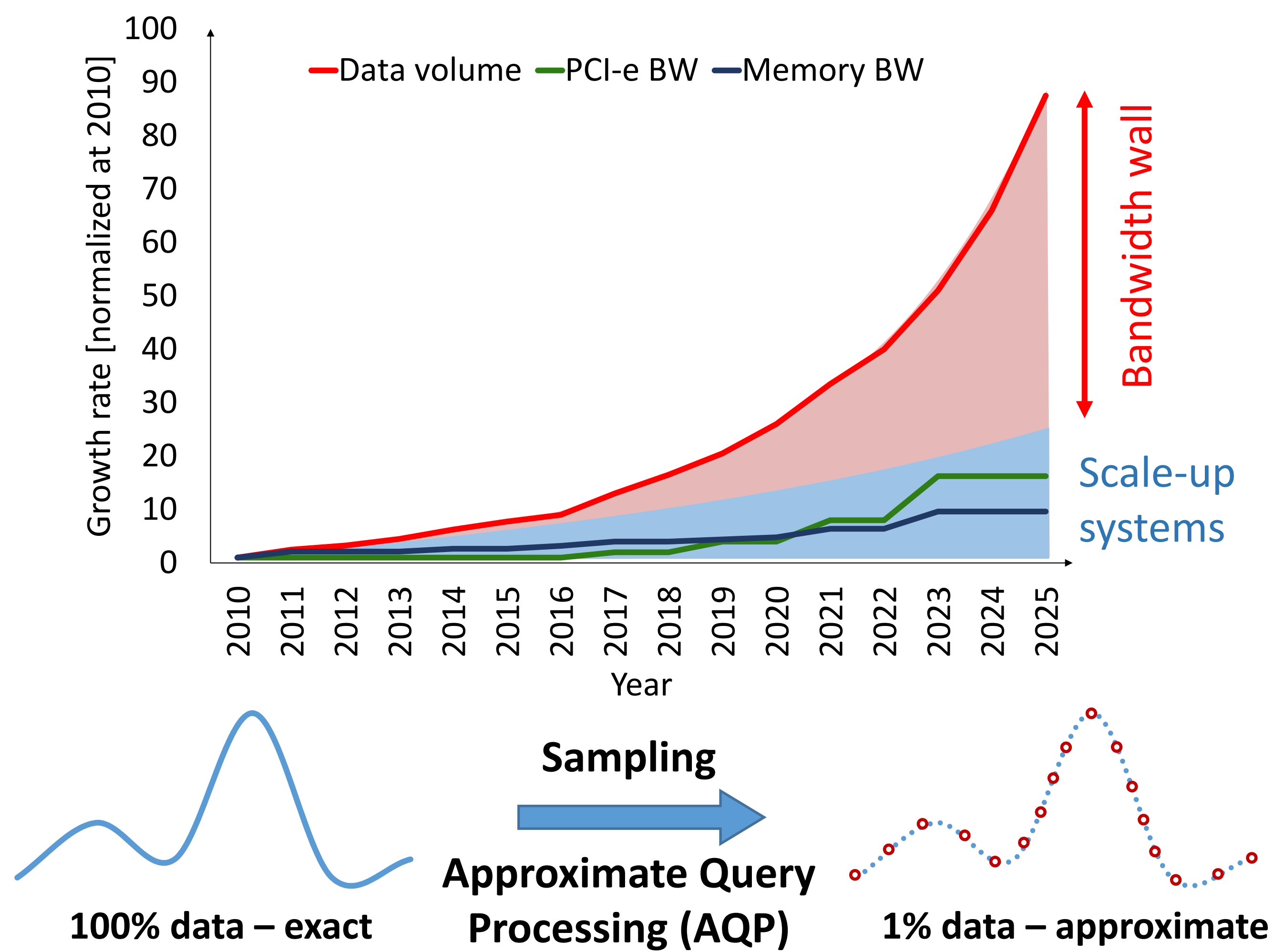


# Sampling-Based AQP in Modern Analytical Engines

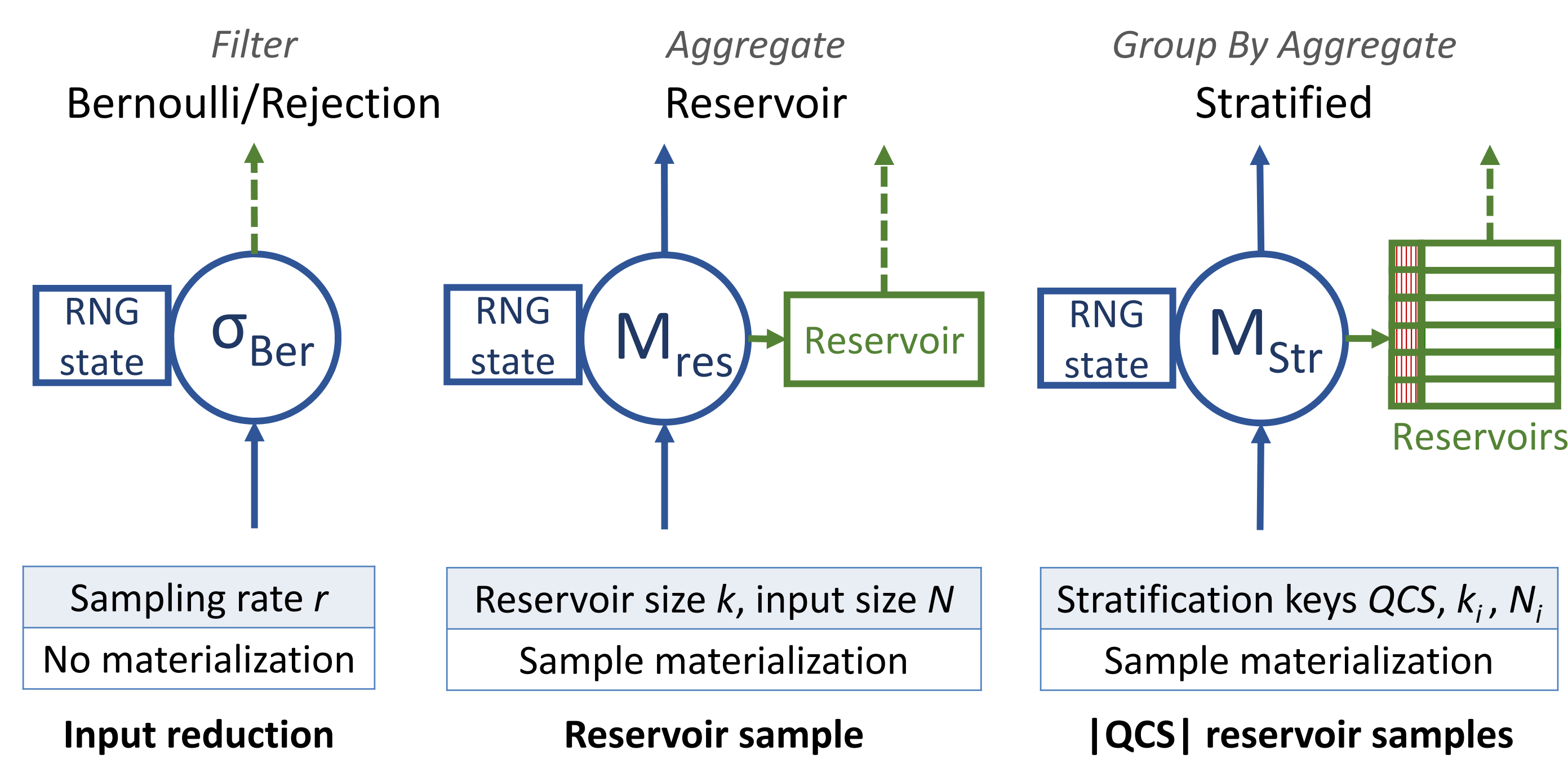
Viktor Sanca and Anastasia Ailamaki

## Interactive analytics is an elusive goal



AQP + modern systems = faster analytics

## Sampling operator design



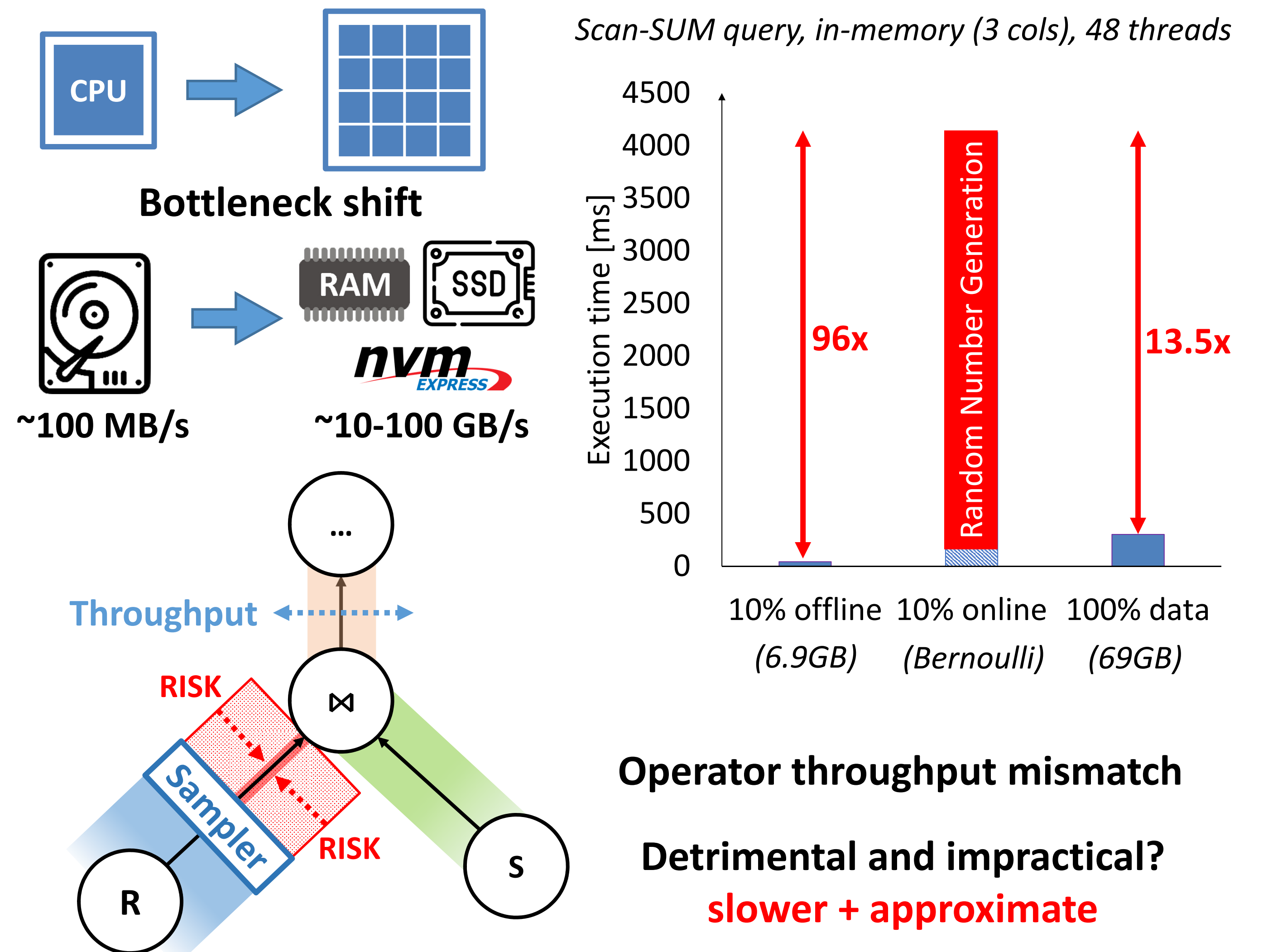
Seamless integration with the existing system

Just-in-time sampling triggered by query operators

Match the throughput of corresponding relational operators

Goal: low overhead side-effect of query execution

## Sampling in the critical path of execution

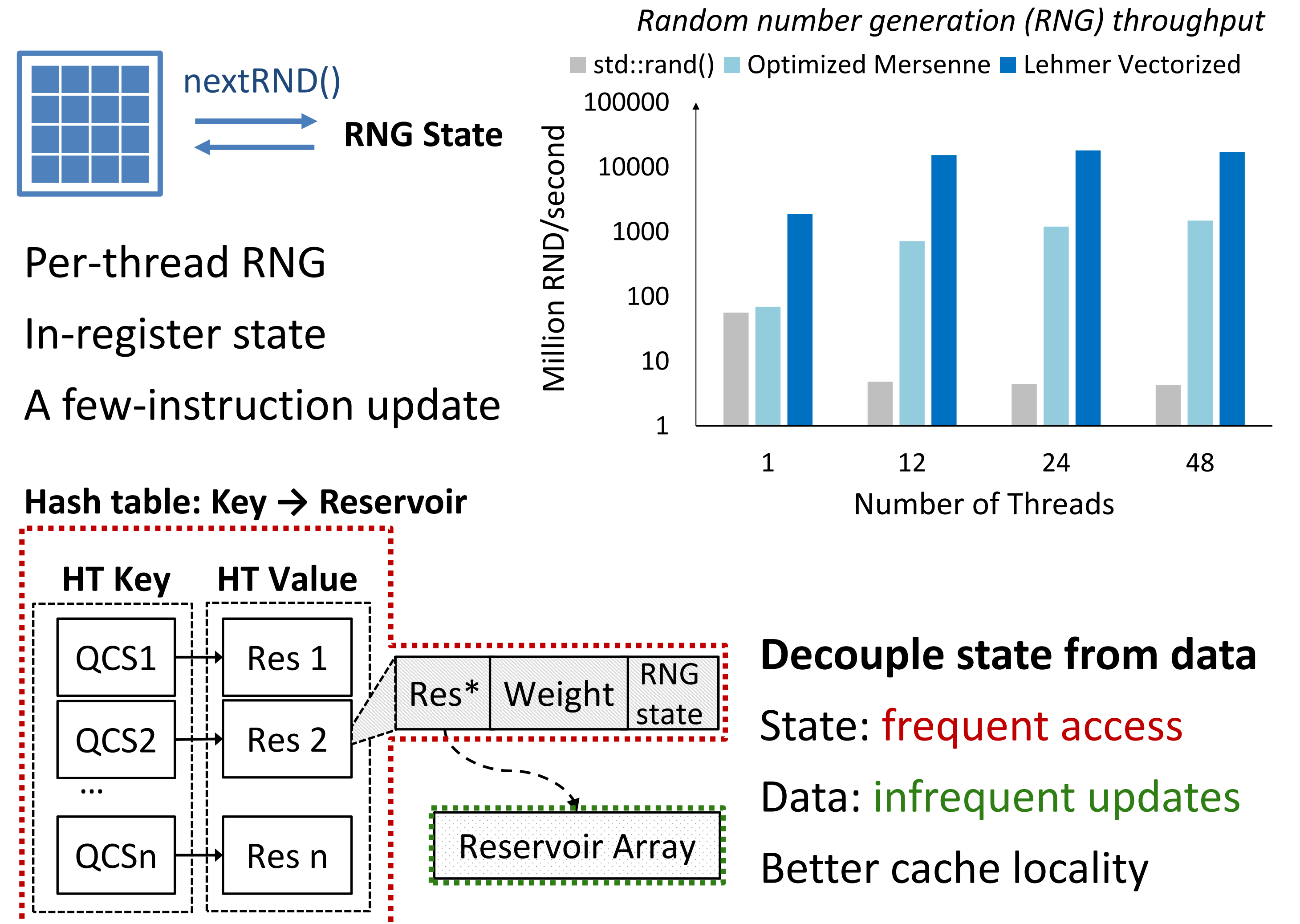


Operator throughput mismatch

Detrimental and impractical?  
slower + approximate

Modern analytics require hardware-consciousness

## Common operations and access patterns



Decouple state from data

State: frequent access

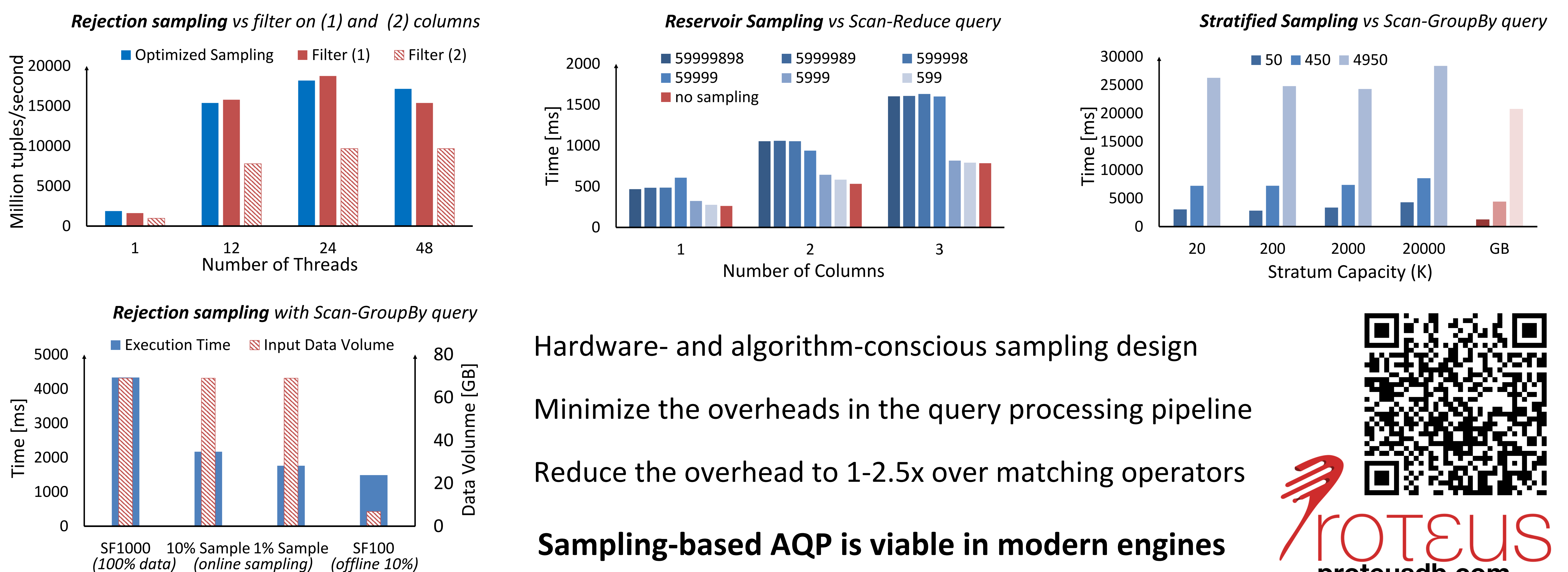
Data: infrequent updates

Better cache locality

Co-design the physical operators with algorithms

## Exposing the bottlenecks: sampling inside an in-memory scale-up analytical engine

Setup: dual socket Intel Xeon Gold 5118 (2x12 cores), 384GB RAM Data: SSB with 600M (SF100) and 6B (SF1000) tuples in fact table, 1 binary column has ~2.3/23GB



Hardware- and algorithm-conscious sampling design

Minimize the overheads in the query processing pipeline

Reduce the overhead to 1-2.5x over matching operators

Sampling-based AQP is viable in modern engines

