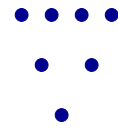
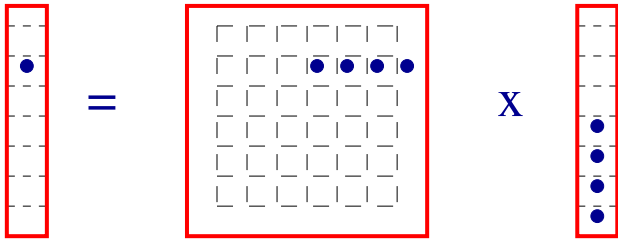


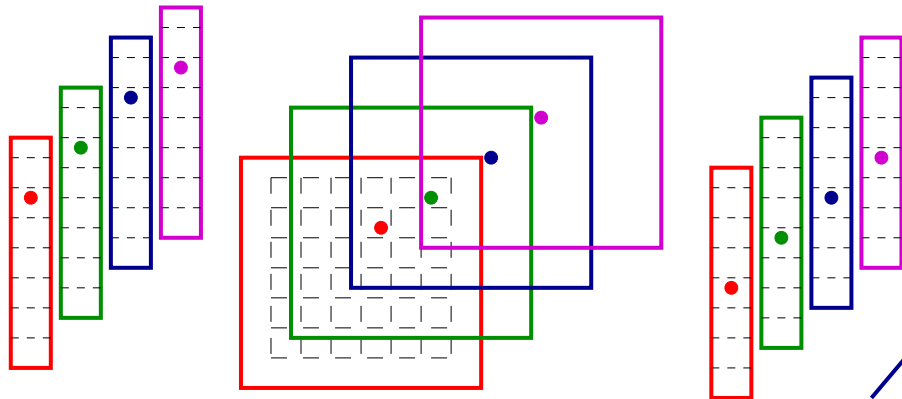
Unifying CPU and GPU programming models

Vector = Matrix x Vector



Reduction of vector sum
is bottleneck for small N

Many vectors = many matrices x many vectors



No reduction or SIMD lane
crossing operations.

SIMD interleave

