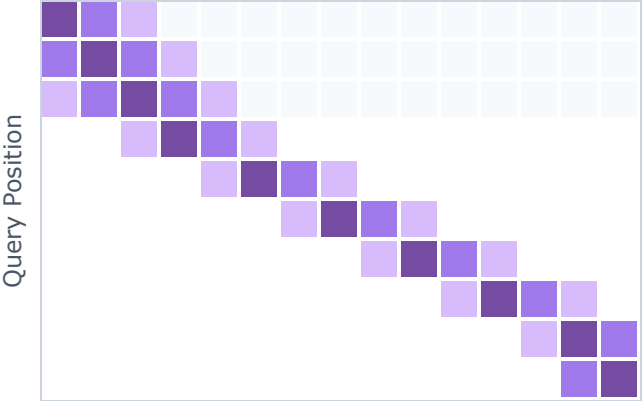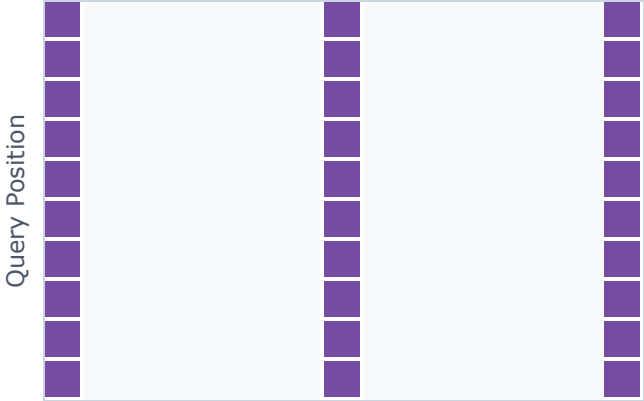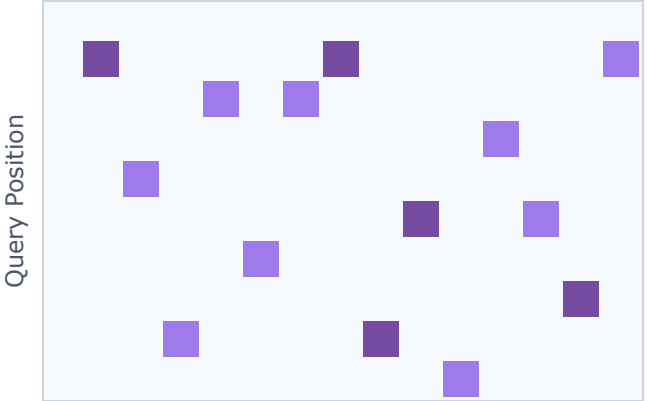# Learned Attention Patterns

## Local Attention



Key Position

Attends to nearby tokens
~60% of heads

## Positional Attention



Key Position

Attends to specific positions
~20% of heads

## Semantic Attention



Key Position

Attends to related content
~20% of heads

## Attention Weight

Low (0.0)    0.3    0.6    High (1.0)

## Optimization Opportunity

Local patterns enable sparse attention: skip low-weight computations
50-80% computation reduction with <1% accuracy impact