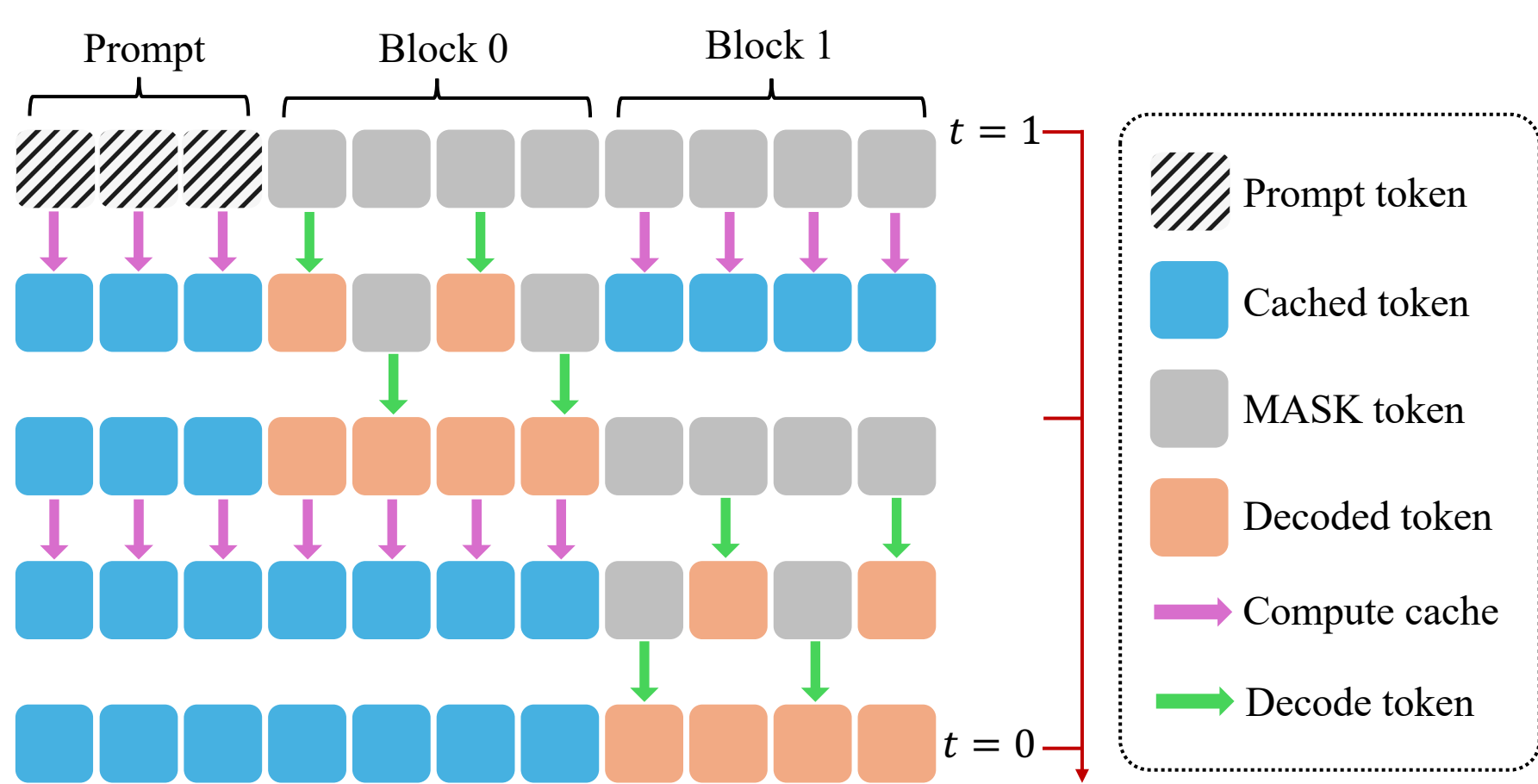(a) Prefix KV Cache for block-wise generation.

(b) **DualCache**: Bidirectional KV cache contains prefix and suffix Cache.

Prompt token

Cached token

MASK token

Decoded token

Compute cache

Decode token