# Flex Attention

Authors present method combining a lot of different attention implementations. During research, authors found that a many attention variants can be defined as score modification applied on the intermediate score matrix before conducting softmax. Score modification is called *score_mod*, attention mask is called *mask_mod*.

FlexAttention leverages block sparsity and employs a pre-computed BlockMask. BlockMask is a small matrix that tracks block-level sparsity on wether a tiled score matrix block is fully masked out(important thing is that Blockmask is generated not during runtime, but during compilation time). This division in scoremod and maskmod is chosen as scoremod is unified for most of the attentions while maskmod provides extra semantic information. Authors use TorchDynamo firstly to build computation graph, after that Triton templates are used for high perfomance. TorchInductor translates these subgraphs into Triton code, dynamically injecting both forward and backward score modification operations into the predefined templates at runtime. In terms of sparsity authors define 2 types of blocks:
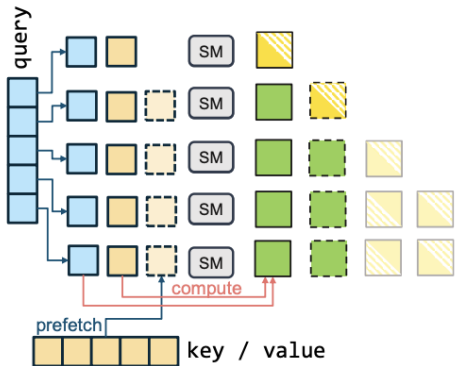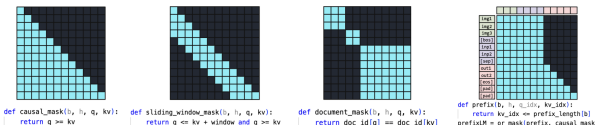
- Partial blocks
- Full blocks





Figure 4. Scheduling full and partial blocks to SM.