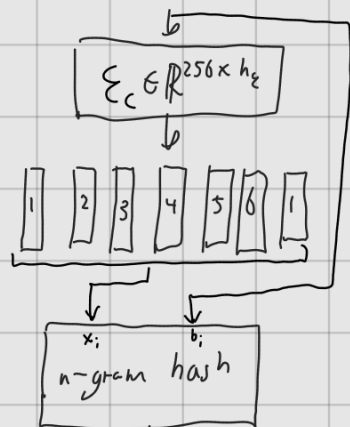


$b \in \mathbb{B}^{n_b \times 256}$

t	h	e	-	c	a	t
0	1	0	1	0	1	0
0	0	1	1	0	0	0
0	0	0	0	1	1	0
1	2	3	4	5	6	1



$$e_i = \frac{1}{|n|+1} \left[x_i + \sum_n E_h^{\text{hash}}(\text{Hash}(g_i, n)) \right]$$

Annotations:
 x_i : i -th byte
 E_h^{hash} : map integer index to vector embedding
 $\text{Hash}(g_i, n)$: map integer sequence to an integer index
 \sum_n : sum over all n n -gram windows

just an integer sequence

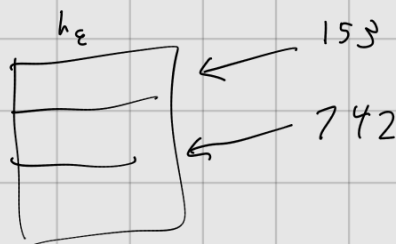
$n = \{3, 4, 5, 6, 7, 8\}$

$g_{i,n} = \{b_{i-n+1}, \dots, b_i\}$

Annotation: $n-1$ previous bytes before the i -th byte and the i -th byte

$\text{Hash}(g_{i,n}) = \text{RollPolyHash}(g_{i,n}) \% |E_n^{\text{hash}}|$

Annotation: integer in $[0, |E_n^{\text{hash}}|]$ the size of the hash table

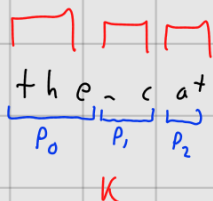
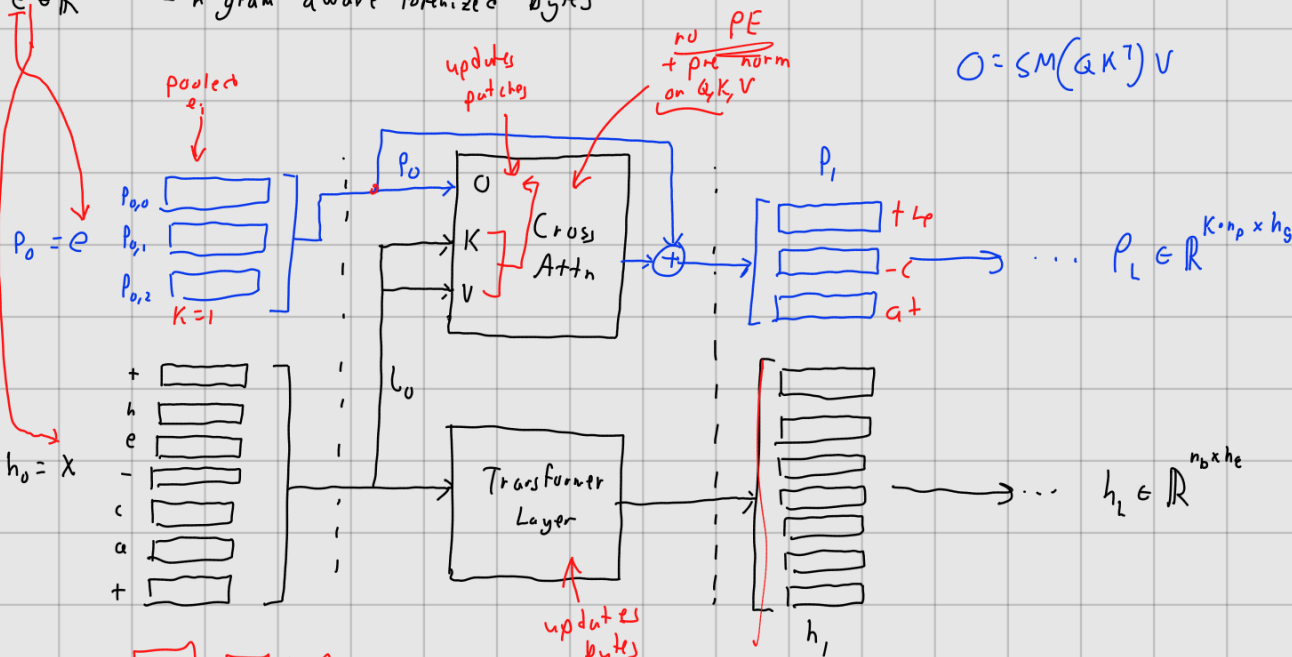


$b \in \mathbb{B}^{n_b \times 256}$ - raw bytes

$x \in \mathbb{R}^{n_b \times h_e}$ - tokenized bytes

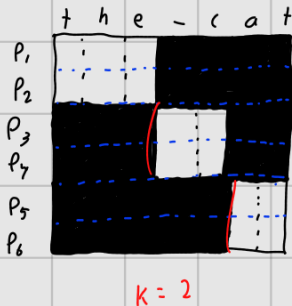
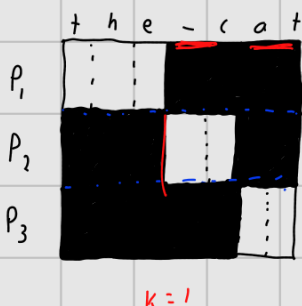
$e \in \mathbb{R}^{n_b \times h_e}$ - n -gram aware tokenized bytes

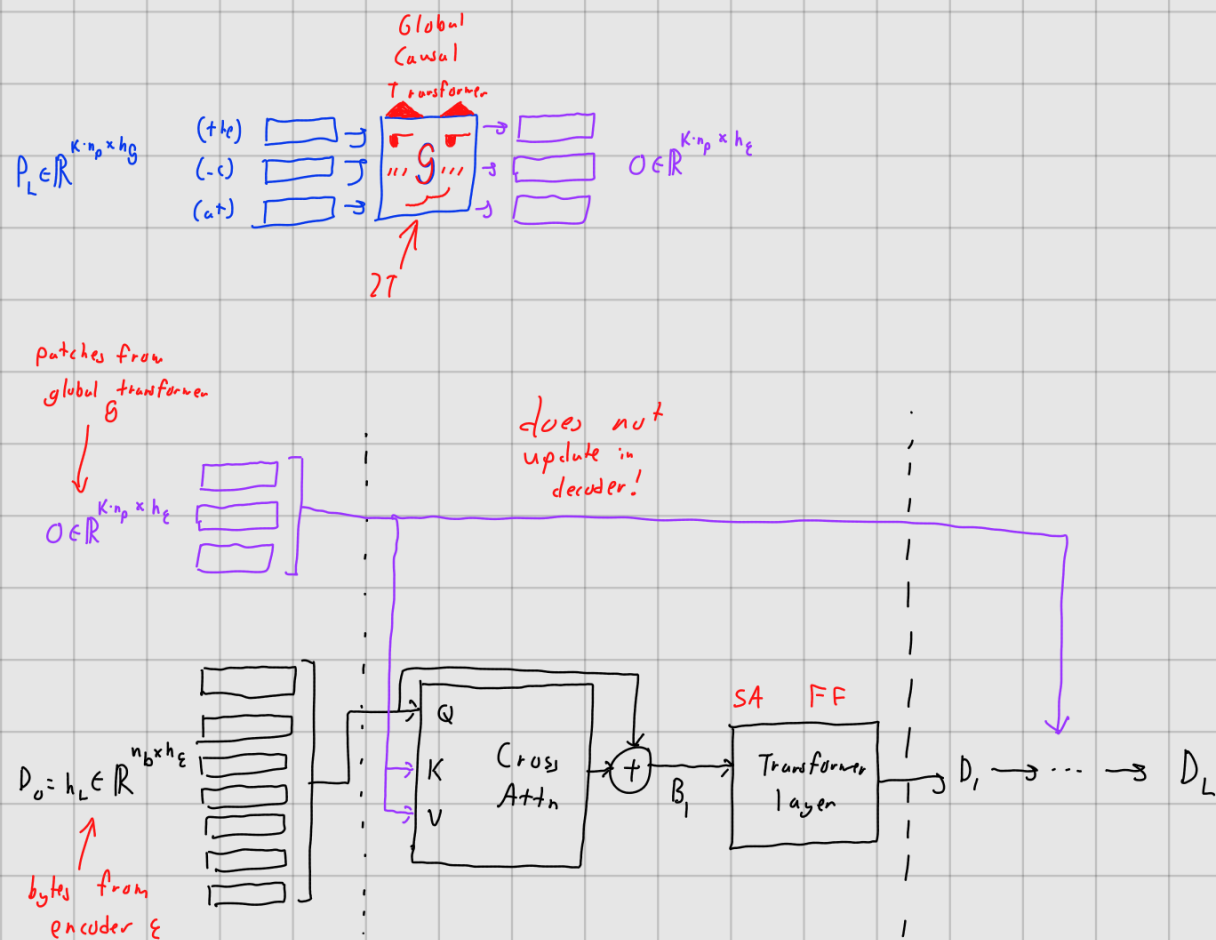
$O = SM(QK^T)V$



bytes send info to patches!

K - # of tokens per patch





Sampling:

```

cur_seq = ''
for byte in num_bytes_to_pred:
    patch_mask = get_patches(cur_seq)
    /, e = encode_and_hash(cur_seq)
    p = pool_patches(e)
    h = X
    for enc_blk in encoder:
        p, h = enc_blk(p, h)
    for global_blk in Global:
        p = global_blk(p)
    O = p
    for dec_blk in decoder:
        h = dec_blk(h, O)
    new_byte = sample(SM(h[-1]))
    cur_seq = cur_seq + new_byte

```

cur probably speed up with KV cache for each byte

only needs to be redone for newest byte

return for last patch

update patches and bytes

update patches

update bytes

add new byte to seq

can probably improve via KV cache. Note the last token cannot be cached until the patch is complete.