Causal Attention with Lookahead Keys

Zhuoqing Song, Peng Sun, Huizhuo Yuan, Quanquan Gu

October XX, 2025

Motivation

- Scaling law runs up against a practical bottleneck high quality data
- How to improve models' performance under limited token budget?

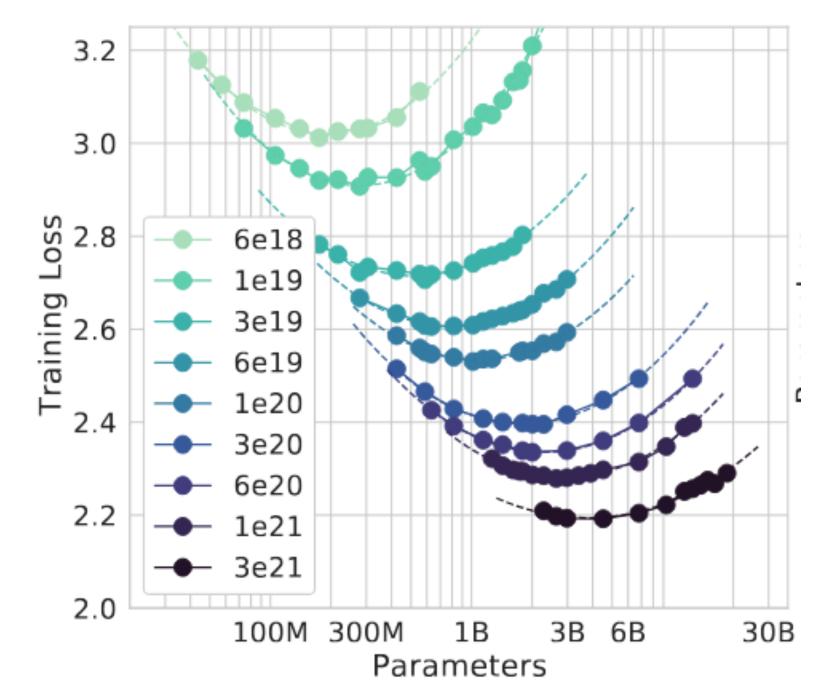


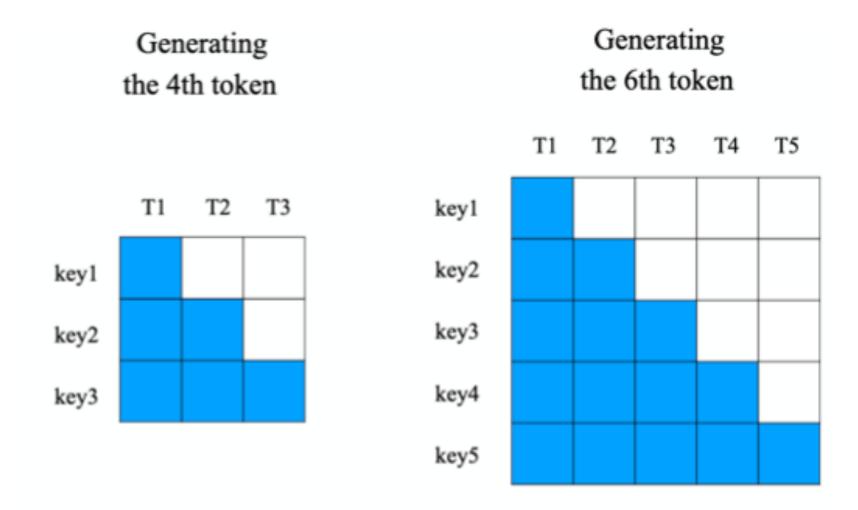
Figure source: https://arxiv.org/pdf/2203.15556

Shortcomings of Standard Causal Attention

Standard causal attention in recurrent form

$$\text{causal-attention}(\boldsymbol{X}^t) = \text{softmax}\left(\frac{\boldsymbol{q}_t \boldsymbol{K}_t^\top}{\sqrt{d}}\right) \boldsymbol{V}_t = \frac{\sum_{s=1}^t \exp\left(\boldsymbol{q}_t \boldsymbol{k}_s^\top/\sqrt{d}\right) \boldsymbol{v}_s}{\sum_{s=1}^t \exp\left(\boldsymbol{q}_t \boldsymbol{k}_s^\top/\sqrt{d}\right)} \in \mathbb{R}^{1 \times d}.$$

Receptive fields of keys in standard causal attention



Standard Causal Attention in Recurrent Form

• Example: "The horse raced past the barn fell."

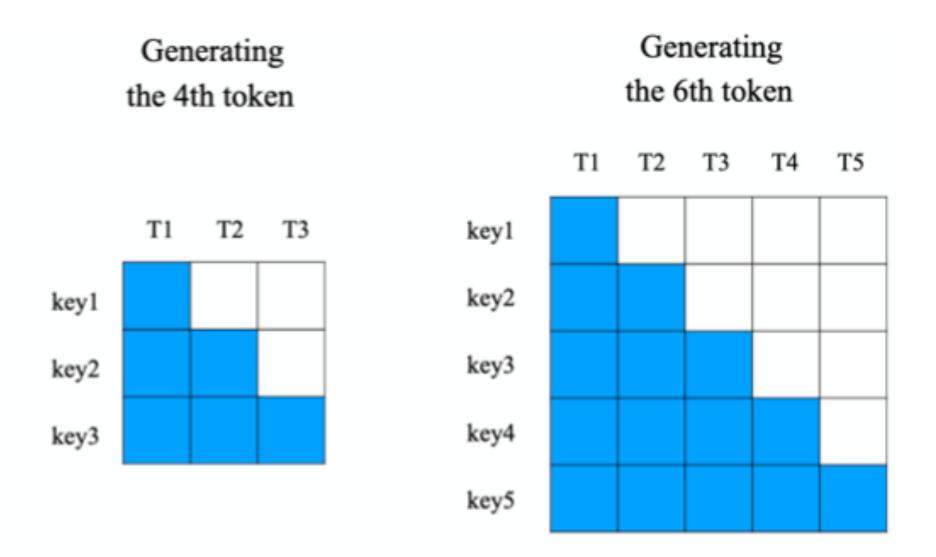
Input

 q_t past $q_t K_t^{ op}$ softmax $\left(\frac{q_t K_t^{ op}}{\sqrt{d}}\right)$ The K_t horse raced softmax $\left(\frac{q_t K_t^{ op}}{\sqrt{d}}\right) V_t$

Output

Shortcomings of Standard Causal Attention

Receptive fields of keys in standard causal attention



- Causal mask blocks each token's access to its future information
 - Hurt natural language understanding (BERT vs. GPT)
 - Noise from precedent context
 - (Indirect evidence) diffusion LLMs are using bi-directional attention

Examples: Shortcomings of Causal Masking

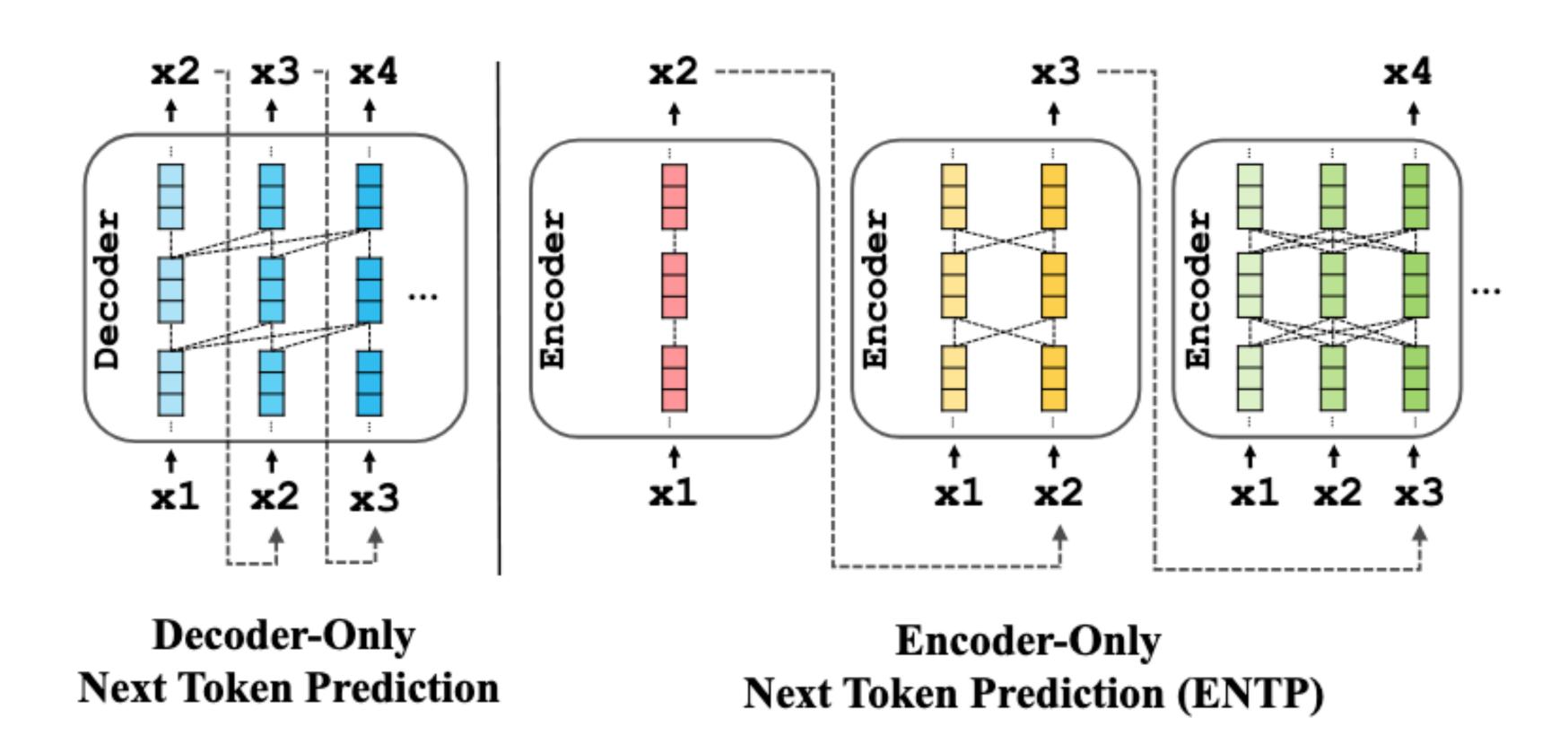
- Garden-path sentence
 - The old man the boat.
 - The complex houses married and single soldiers and their families.
- Question/focus at the end of inputs
- Variable assignment problem
 - x = 1; y = 2; x = 3; y = 5; x = ?

Related work: Sentence Embedding

- Backward Dependency Enhanced LLM
- Echo Embeddings
- Re-Reading

Related work: Encoder-only Next Token Prediction (ENTP)

• NTP by encoder-only Transformers (on KV cache, cubic training complexity)

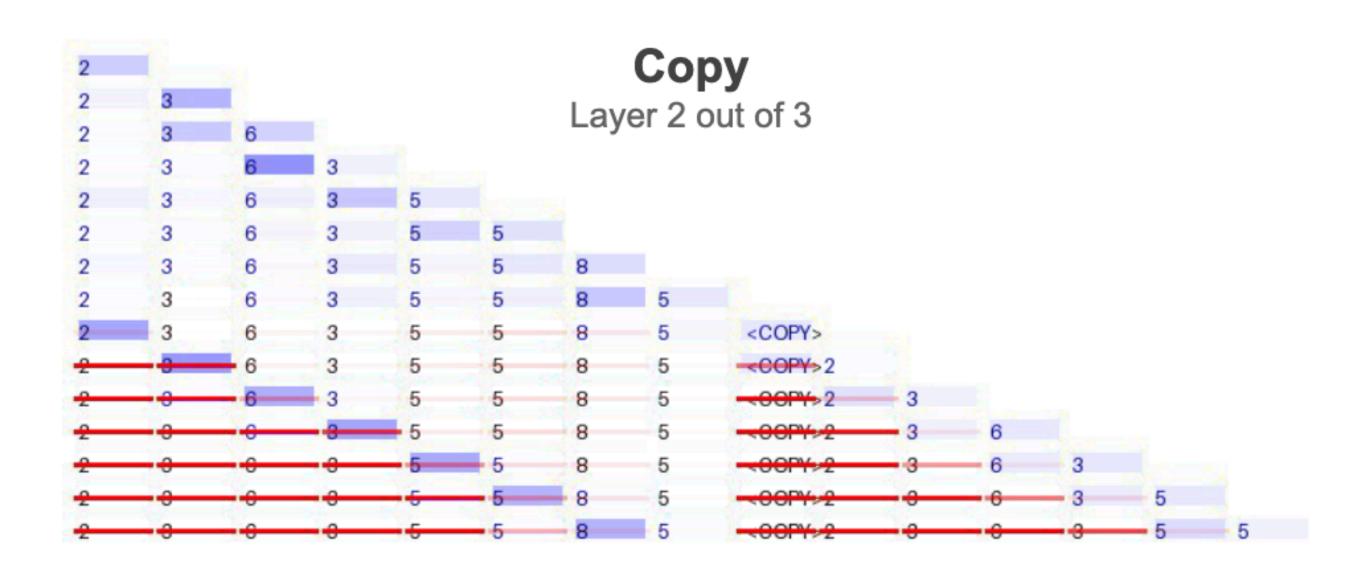


Related work: Selective Attention

Modify attention matrix, save KV cache

SelectiveAttention
$$(Q, K, V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}} - F)V$$

$$F = Accumulate(Constrain(S))$$



Related work: PaTH Attention

Data-dependent positional encodings

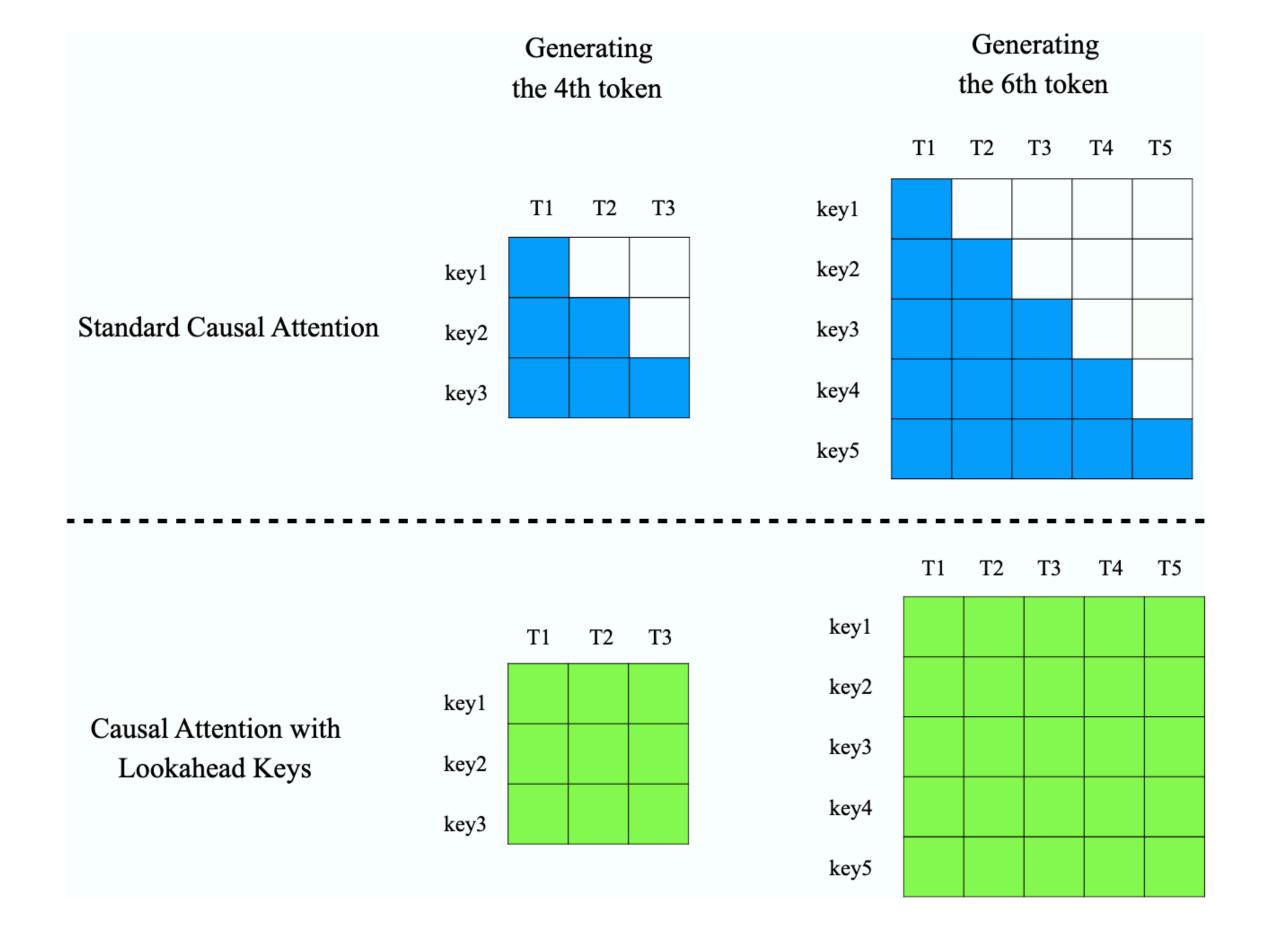
$$\mathbf{A}_{ij} \propto \exp\Bigl(\mathbf{k}_j^{ op}\Bigl(\prod_{s=j+1}^i \mathbf{H}_s\Bigr)\mathbf{q}_i\Bigr),$$
 $\mathbf{H}_t = \mathbf{I} - eta_t \mathbf{w}_t \mathbf{w}_t^T$

Inference of PaTH

$$\mathbf{k}_i^{(t)} \leftarrow (\mathbf{I} - \beta_t \mathbf{w}_t \mathbf{w}_t^{\top}) \mathbf{k}_i^{(t-1)}$$
 for all $i < t$,

Our Approach: Expanding Receptive Fields

• Should we expand receptive fields of Q, K or V?



Our Approach: Expanding Receptive Fields

• Should we expand receptive fields of Q, K or V?

The same reason with using KV cache instead of Q cache!

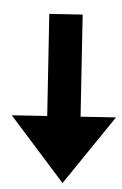
Why are we updating keys rather than values?

Updating keys has an equivalent form which enables efficient parallel training.

If attentions scores are accurate `enough', information mix in values is unnecessary.

Our Approach: Expanding Receptive Fields

- Expanding receptive fields of keys
- Maintaining efficient training and inference



- Update keys during inference
- Without materializing keys at each position during training

CASTLE in Recurrent Form (Overview)

CASTLE

Causal-key attention score

$$oldsymbol{s}_t^C = rac{oldsymbol{q}_t^C oldsymbol{K}_t^{C^ op}}{\sqrt{d}} \in \mathbb{R}^{1 imes t}$$

Lookahead-key attention score

$$oldsymbol{s}_t^U = rac{oldsymbol{q}_t^C oldsymbol{U}^{t^ op}}{\sqrt{d}} \in \mathbb{R}^{1 imes t}$$

Attention weights

$$oldsymbol{p}_t = \operatorname{softmax}\left(oldsymbol{s}_t^C - \operatorname{SiLU}\left(oldsymbol{s}_t^U
ight)\right) \in \mathbb{R}^{1 imes t}$$

Outputs

attention
$$(\boldsymbol{X}^t) = \boldsymbol{p}_t \boldsymbol{V}_t^C \in \mathbb{R}^{1 \times d}$$
.

Standard Causal Attention

Causal-key attention score

$$oldsymbol{s}_t = rac{oldsymbol{q}_t oldsymbol{K}_t^ op}{\sqrt{d}} \in \mathbb{R}^{1 imes t}$$

Attention weights

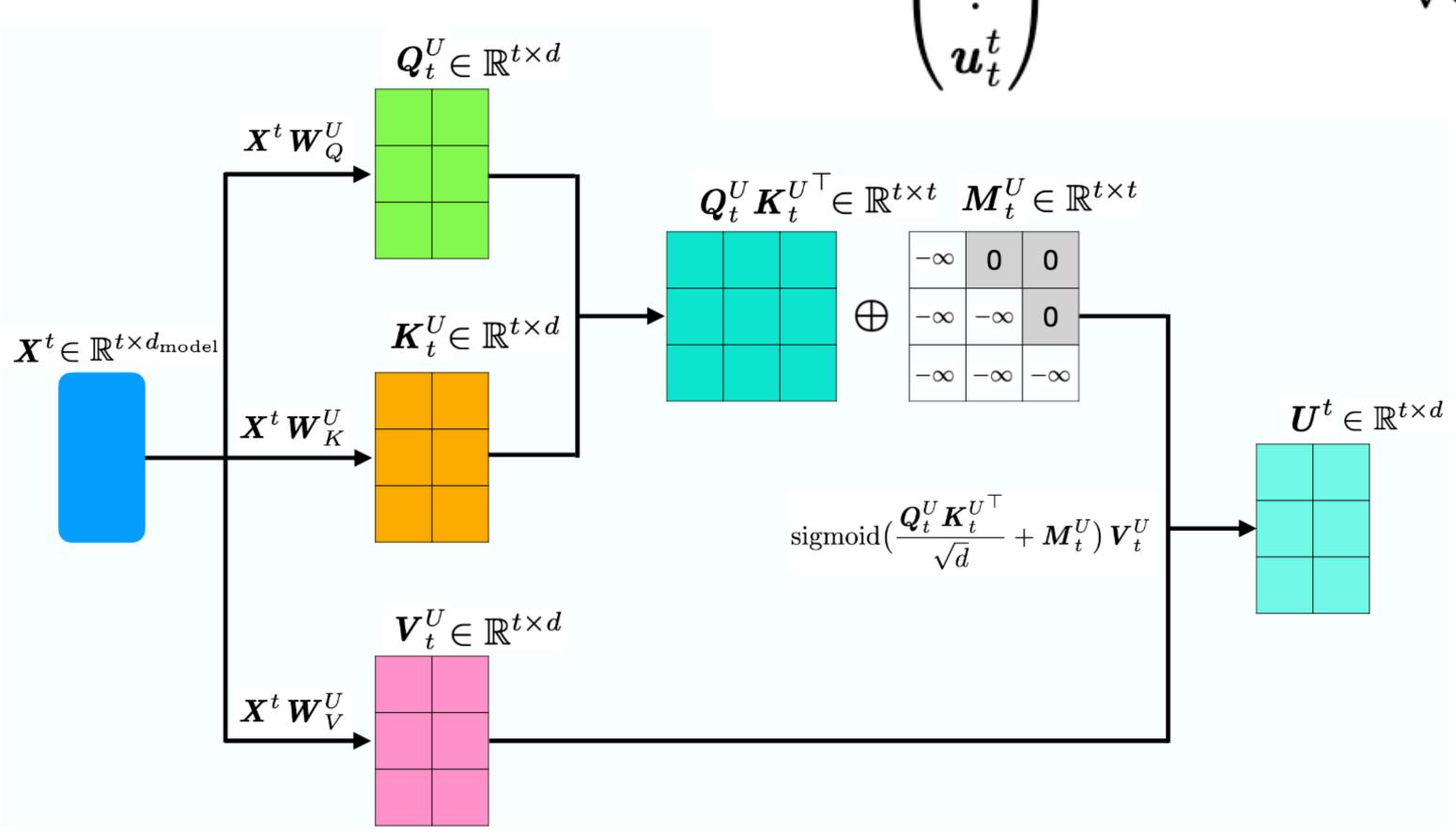
$$\boldsymbol{p}_t = \operatorname{softmax}\left(\boldsymbol{s}_t\right) \in \mathbb{R}^{1 \times t}$$

Outputs

attention
$$(\boldsymbol{X}^t) = \boldsymbol{p}_t \, \boldsymbol{V}_t \in \mathbb{R}^{1 \times d}$$
.

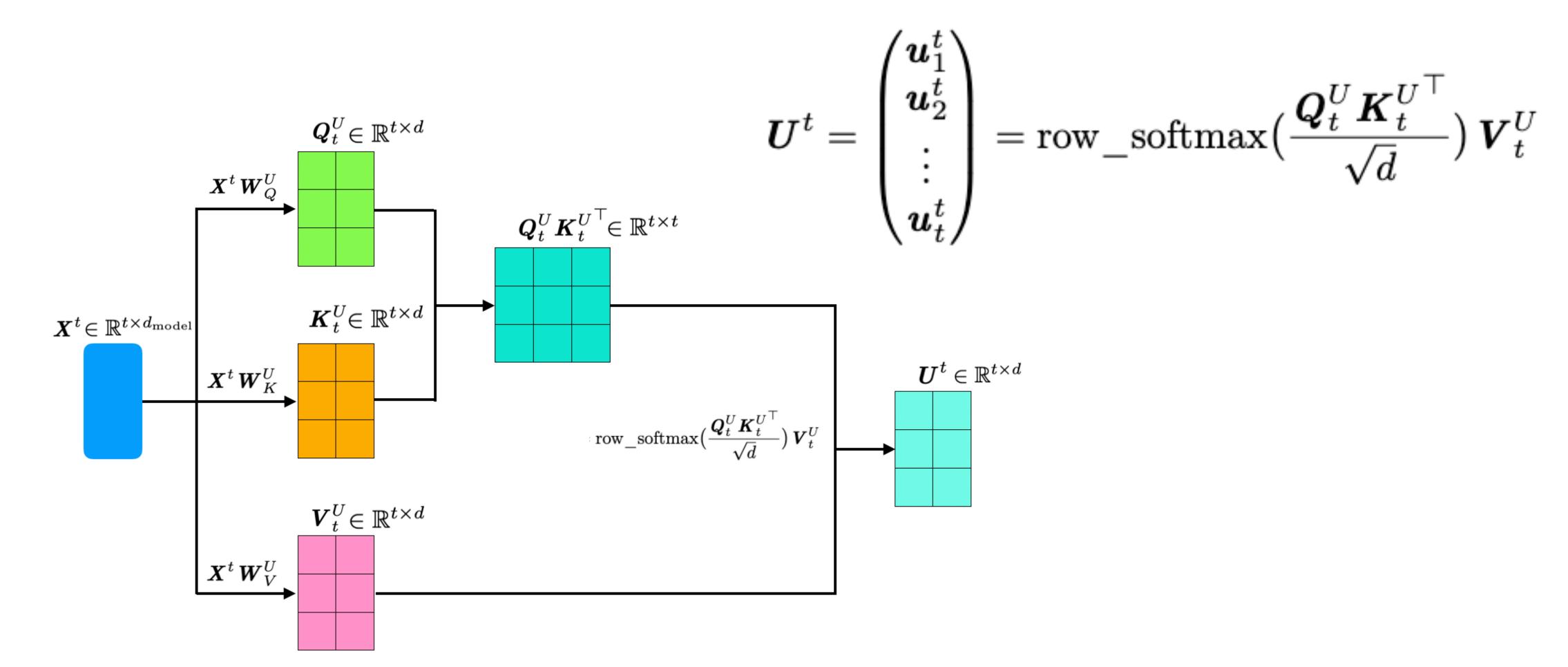
Formal Definition of Lookahead Keys

• Similar to attention mechanism
$$egin{aligned} oldsymbol{U}^t = \begin{pmatrix} oldsymbol{u}_1^t \ oldsymbol{u}_2^t \ \vdots \ oldsymbol{u}_t^t \end{pmatrix} = \operatorname{sigmoid}ig(rac{oldsymbol{Q}_t^U oldsymbol{K}_t^{U^\top}}{\sqrt{d}} + oldsymbol{M}_t^U ig) oldsymbol{V}_t^U \in \mathbb{R}^{t imes d}, \end{aligned}$$



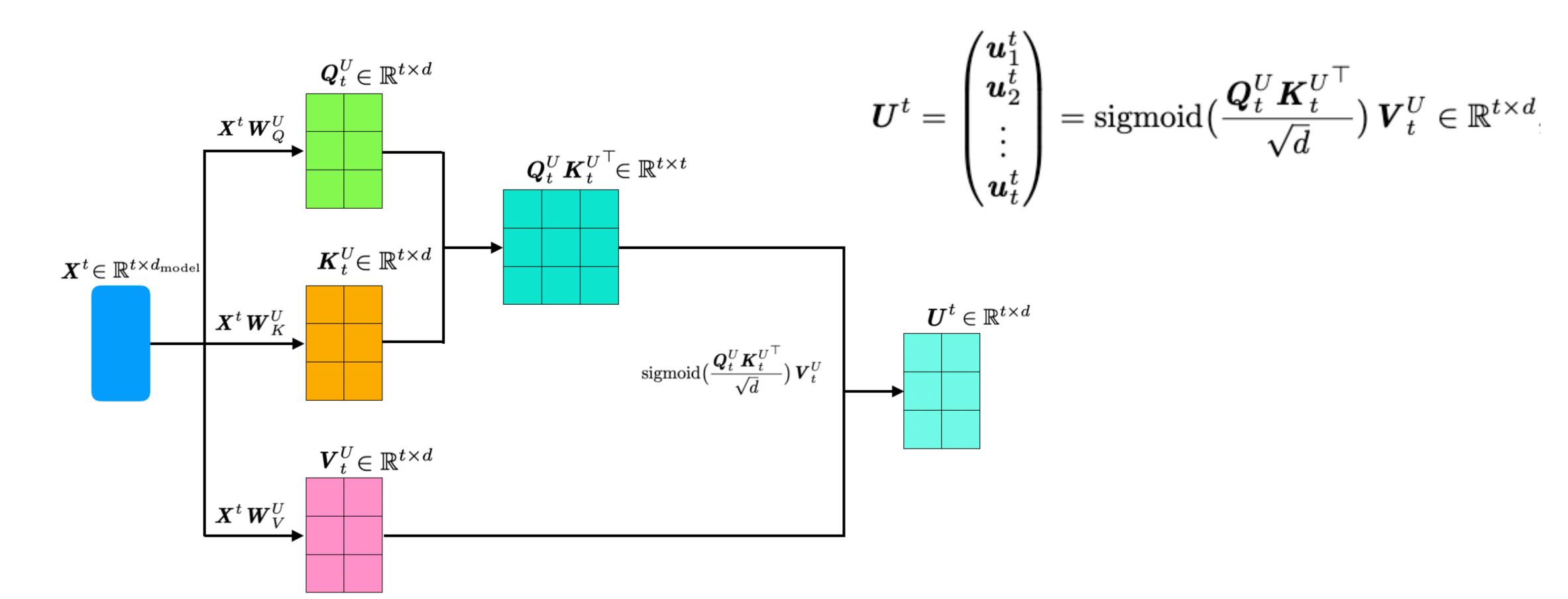
Failed Attempt of Lookahead keys

• Encode all information from token 1 to token t: numerically unstable



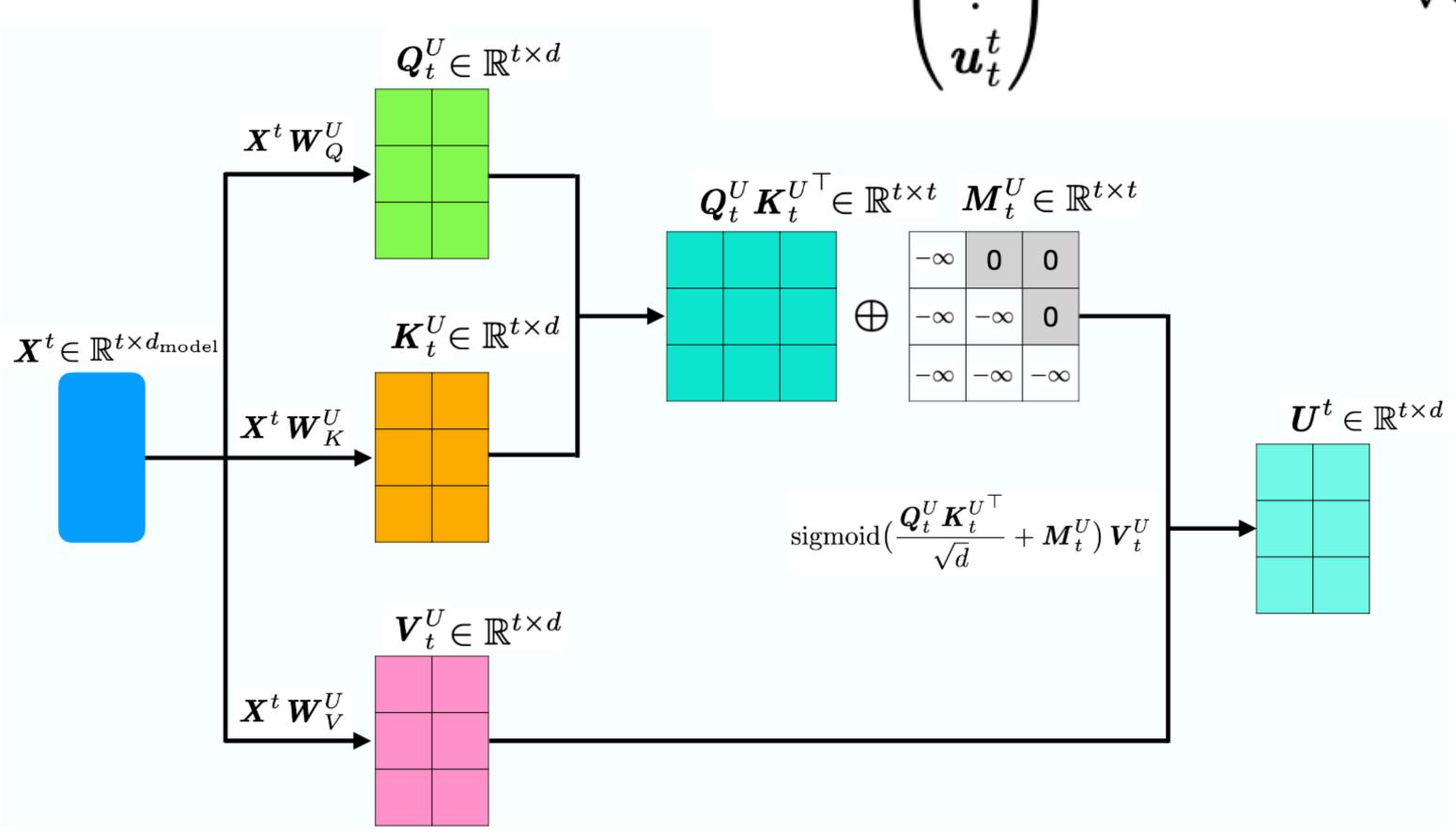
Failed Attempt of Lookahead keys

• Replace softmax by sigmoid: training instability (blow up easily)



Formal Definition of Lookahead Keys

• Similar to attention mechanism
$$egin{aligned} oldsymbol{U}^t = \begin{pmatrix} oldsymbol{u}_1^t \ oldsymbol{u}_2^t \ \vdots \ oldsymbol{u}_t^t \end{pmatrix} = \operatorname{sigmoid}ig(rac{oldsymbol{Q}_t^U oldsymbol{K}_t^{U^\top}}{\sqrt{d}} + oldsymbol{M}_t^U ig) oldsymbol{V}_t^U \in \mathbb{R}^{t imes d}, \end{aligned}$$



Definition of Causal Keys

Causal keys are the same as the keys in standard causal attention

$$m{K}_t^C = egin{pmatrix} m{k}_1^C \ m{k}_2^C \ dots \ m{k}_t^C \end{pmatrix} = m{X}^t \, m{W}_K^C \in \mathbb{R}^{t imes d}$$

Ablation Studies on Causal Keys

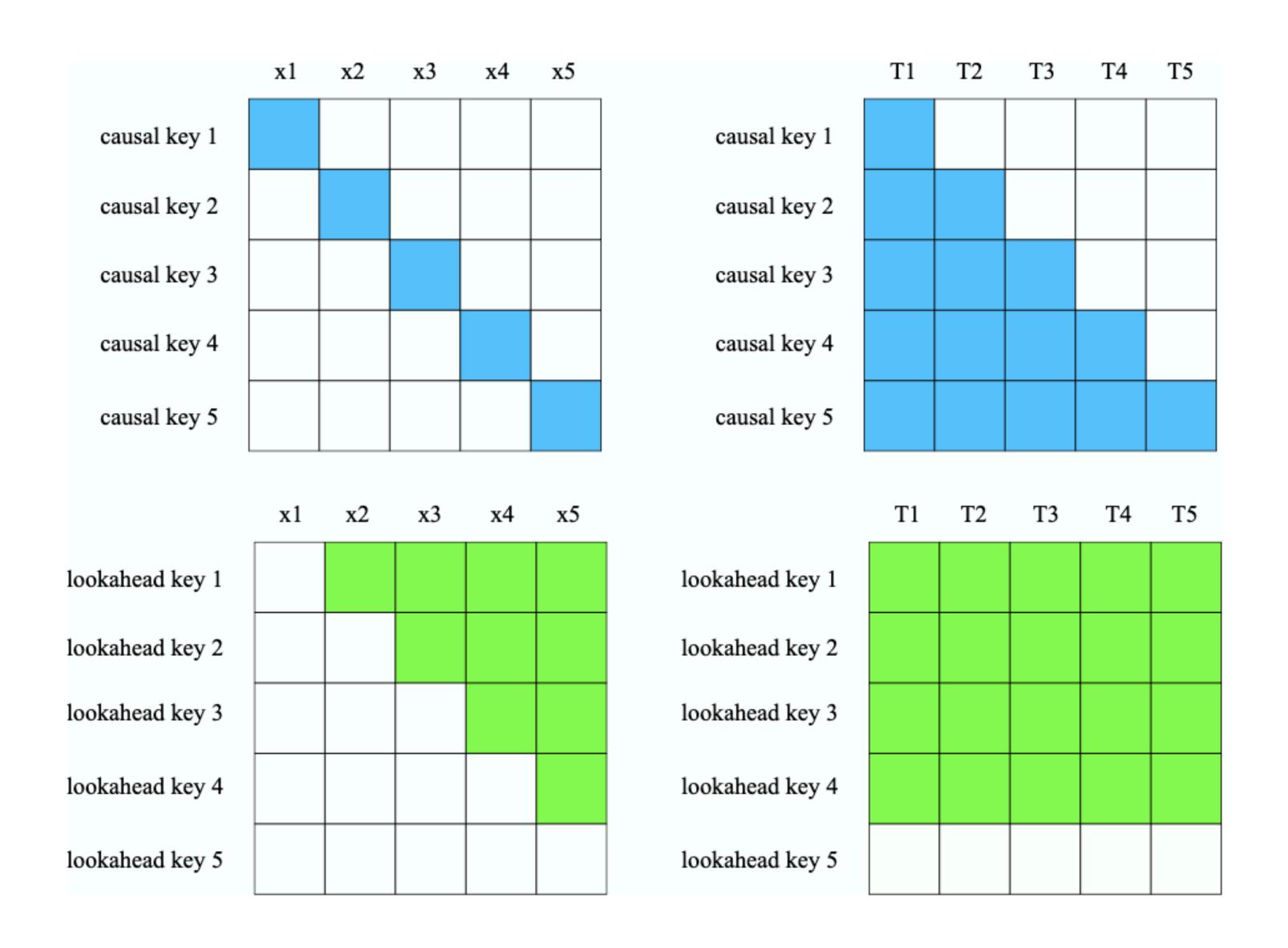
- Whether we need causal keys?
 - Aligning parameters

CASTLE TYPE	$n_{ m params}$	$n_{ m layers}$	$d_{ m model}$	$n_{ m heads}$	d
CASTLE	120M	12	768	6	64
CASTLE w/o causal keys	120M	12	768	7	64

Training & validation loss

	Tr	rain	Eval		
	Loss	PPL	Loss	PPL	
CASTLE	2.913	18.417	2.920	18.541	
${\it CASTLE~w/o~causal~keys}$	3.006	20.213	3.021	20.505	

Comparison between Causal Keys and Lookahead Keys



CASTLE in Recurrent Form

Causal-key attention score

$$oldsymbol{s}_t^C = rac{oldsymbol{q}_t^C oldsymbol{K}_t^{C^ op}}{\sqrt{d}} \in \mathbb{R}^{1 imes t}$$

Lookahead-key attention score

$$oldsymbol{s}_t^U = rac{oldsymbol{q}_t^C oldsymbol{U}^{t^{-1}}}{\sqrt{d}} \in \mathbb{R}^{1 imes t}$$

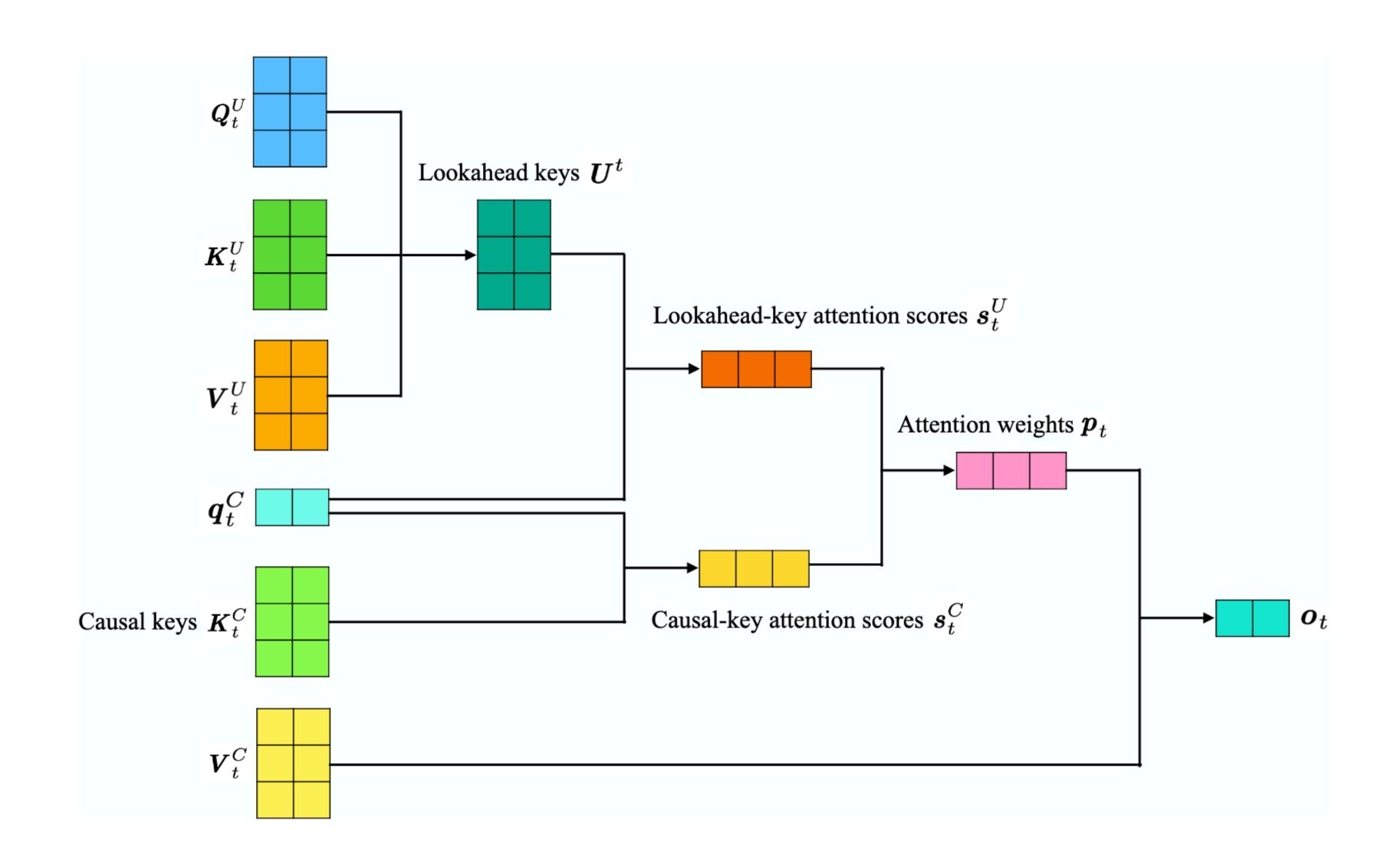
Attention weights

$$\boldsymbol{p}_{t} = \operatorname{softmax}\left(\boldsymbol{s}_{t}^{C} - \operatorname{SiLU}\left(\boldsymbol{s}_{t}^{U}\right)\right) \in \mathbb{R}^{1 \times t}$$

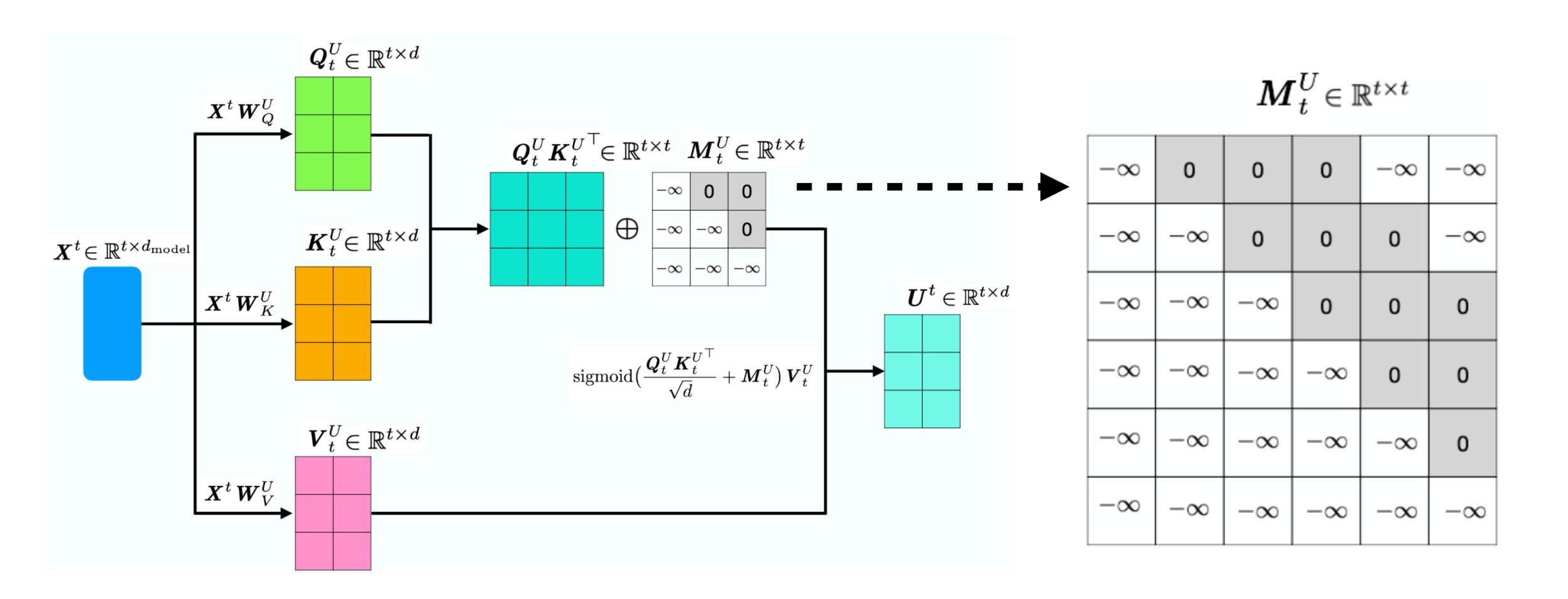
outputs

attention
$$(\boldsymbol{X}^t) = \boldsymbol{p}_t \boldsymbol{V}_t^C \in \mathbb{R}^{1 \times d}$$
.

Illustration for CASTLE in Recurrent Form



A variant of CASTLE: CASTLE-SWL



Ablation Studies on Sliding Window Sizes in CASTLE-SWL

• Not sensitive in tested range [128, 1024]

		Tr	rain	Eval		
	$n_{ m params}$	Loss	PPL	Loss	PPL	
Baseline-S	160M	2.892	18.037	2.901	18.197	
CASTLE-SWL128-S	160M	2.883	17.859	2.889	17.971	
CASTLE-SWL256-S	160M	2.888	17.952	2.895	18.092	
CASTLE-SWL512-S	160M	2.885	17.910	2.892	18.023	
CASTLE-SWL1024-S	160M	2.885	17.901	2.892	18.031	

		Th	rain	Eval		
	$n_{ m params}$	Loss	PPL	Loss	PPL	
Baseline-XL	1.310B	2.548	12.779	2.543	12.723	
CASTLE-SWL128-XL	1.304B	2.517	12.393	2.513	12.339	
CASTLE-SWL256-XL	1.304B	2.514	12.353	2.510	12.300	
CASTLE-SWL512-XL	1.304B	2.506	12.255	2.503	12.217	
CASTLE-SWL1024-XL	1.304B	2.514	12.351	2.509	12.294	

$$oldsymbol{M}_t^U \in \mathbb{R}^{t imes t}$$

$-\infty$	0	0	0	$-\infty$	$-\infty$
$-\infty$	$-\infty$	0	0	0	-8
$-\infty$	$-\infty$	$-\infty$	0	0	0
$-\infty$	$-\infty$	$-\infty$	$-\infty$	0	0
$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	0
$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$	$-\infty$

Ablation Studies on SiLU in softmax

- What roles does SiLU play in the combining step? Training stability
- Removing SiLU will cause blowing up when training XL models for both CASTLE and CASTLE-SWL
- Reducing Ir can stabilize training but degrade performance severely

		Train		Eval	
	Learning Rate	Loss	PPL	Loss	PPL
CASTLE-SWL-XL	2×10^{-4} 1×10^{-4}	2.506	12.255	2.503	12.217
CASTLE-SWL-XL w/o SiLU CASTLE-SWL-XL w/o SiLU	5×10^{-5}	$2.523 \\ 2.571$	12.468 13.084	$2.520 \\ 2.571$	$12.424 \\ 13.075$

Revisit the definition of lookahead keys

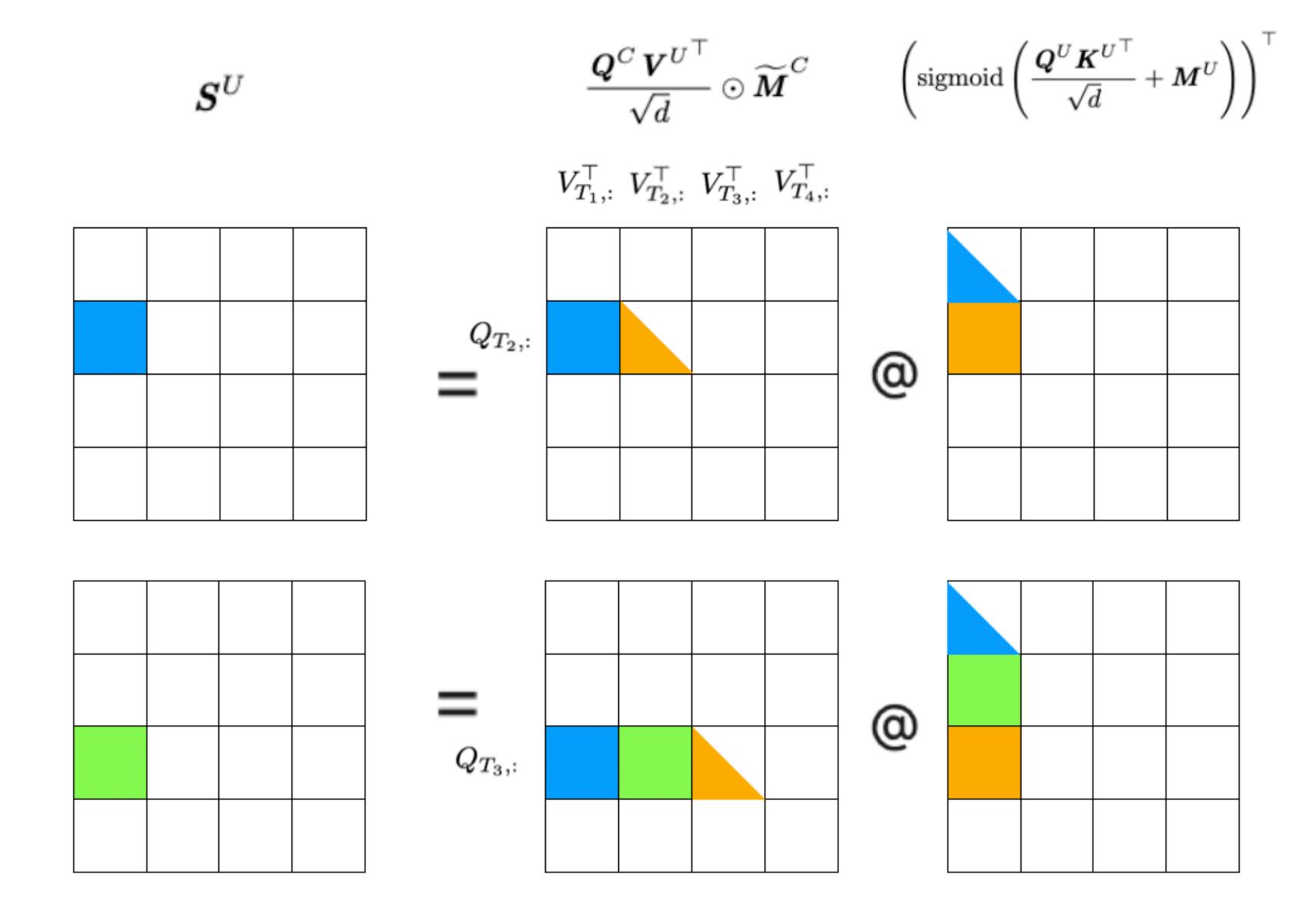
$$oldsymbol{U}^t = egin{pmatrix} oldsymbol{u}_1^t \ oldsymbol{u}_2^t \ oldsymbol{u}_2^t \end{pmatrix} = \operatorname{sigmoid}(oldsymbol{Q}_t^U oldsymbol{K}_t^{U^ op} + oldsymbol{M}_t^U) oldsymbol{V}_t^U \in \mathbb{R}^{t imes d}, \quad oldsymbol{s}_t^U = oldsymbol{q}_t^C oldsymbol{U}^{t^ op} \ oldsymbol{q}_t^{T} \in \mathbb{R}^{1 imes t}, \quad oldsymbol{p}_t = \operatorname{softmax}\left(oldsymbol{s}_t^C - \operatorname{SiLU}\left(oldsymbol{s}_t^U
ight)
ight) \in \mathbb{R}^{1 imes t}, \quad oldsymbol{s}_t^U = oldsymbol{q}_t^U oldsymbol{V}_t^U = oldsymbol{q}_t^U oldsymbol{V}_t^U \in \mathbb{R}^{1 imes t}, \quad oldsymbol{s}_t^U = oldsymbol{q}_t^U oldsymbol{V}_t^U = oldsymbol{s}_t^U oldsymbol{s}_t^U = oldsymbol{q}_t^U oldsymbol{V}_t^U \in \mathbb{R}^{1 imes t}, \quad oldsymbol{s}_t^U = oldsymbol{s}_t^U oldsymbol{s}_t^U = oldsymbol{s}_t^U oldsymbol{s}_t^U oldsymbol{s}_t^U = oldsymbol{s}_t^U oldsymbol{s}_t^U oldsymbol{s}_t^U oldsymbol{s}_t^U = oldsymbol{s}_t^U oldsymbol{s}_t^U oldsymbol{s}_t^U oldsymbol{s}_t^U = oldsymbol{s}_t^U oldsymb$$

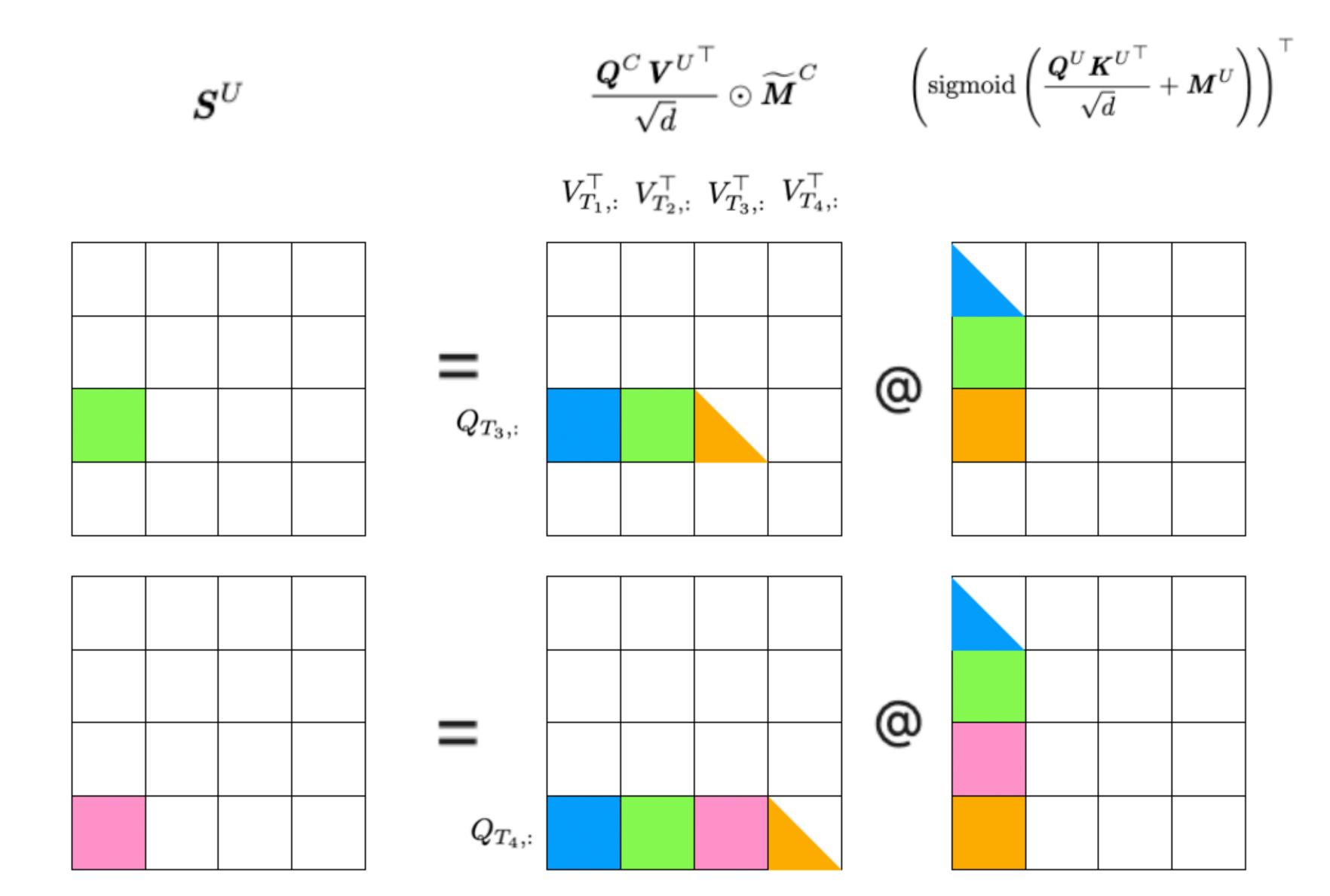
- Equivalent mathematical formulation
 - Lookahead-keys attention scores

$$m{S}^U = \left(rac{m{Q}^Cm{V}^{U^ op}}{\sqrt{d}}\odot\widetilde{m{M}}^C
ight) \left(ext{sigmoid}\left(rac{m{Q}^Um{K}^{U^ op}}{\sqrt{d}}+m{M}^U
ight)
ight)^ op$$

Parallel form of attention outputs

$$Attention(m{X}^L) = row_softmax \left(rac{m{Q}^Cm{K}^C}{\sqrt{d}}^{ op} + m{M}^C - \mathrm{SiLU}(m{S}^U)
ight) m{V}^C$$





$$\begin{aligned} \text{Lookahead-key attention scores} \qquad & \boldsymbol{S}_{T_{j+k},T_{j}}^{U} = \boldsymbol{Q}_{T_{j+k},:}^{C} \sum_{i=j}^{j+k-1} \frac{\boldsymbol{V}_{T_{i},:}^{U}}{\sqrt{d}} \left(\operatorname{sigmoid} \left(\frac{\boldsymbol{Q}_{T_{j},:}^{U} \boldsymbol{K}_{T_{i},:}^{U}}{\sqrt{d}} + \boldsymbol{M}_{T_{j},T_{i}}^{U} \right) \right)^{\top} \\ & + \left(\frac{\boldsymbol{Q}_{T_{j+k},:}^{C} \boldsymbol{V}_{T_{j+k},:}^{U}}{\sqrt{d}} \odot \widetilde{\boldsymbol{M}}_{T_{j+k},T_{j+k}}^{C} \right) \left(\operatorname{sigmoid} \left(\frac{\boldsymbol{Q}_{T_{j},:}^{U} \boldsymbol{K}_{T_{j+k},:}^{U}}{\sqrt{d}} + \boldsymbol{M}_{T_{j},T_{j+k}}^{U} \right) \right)^{\top} \end{aligned}$$

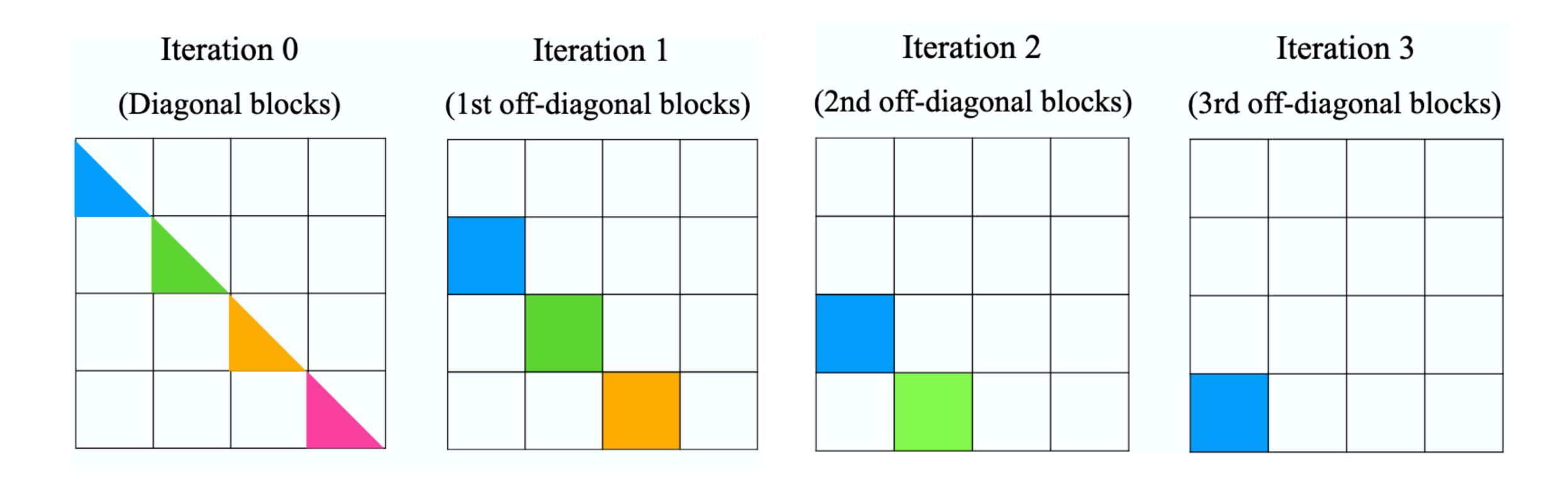
- Utilizing masked low-rank structure
 - Update auxiliary variable

$$\boldsymbol{D}_{:,T_{j}}^{(k)} = \boldsymbol{D}_{:,T_{j}}^{(k-1)} + \boldsymbol{V}_{T_{j+k-1},:}^{U} \top \left(\operatorname{sigmoid} \left(\frac{\boldsymbol{Q}_{T_{j},:}^{U} \boldsymbol{K}_{T_{j+k-1},:}^{U}}{\sqrt{d}} + \boldsymbol{M}_{T_{j},T_{j+k-1}}^{U} \right) \right)^{\top}$$

Compute lookahead-key attention scores

$$\boldsymbol{S}_{T_{j+k},T_{j}}^{U} = \frac{\boldsymbol{Q}_{T_{j+k},:}^{C}\boldsymbol{D}_{:,T_{j}}^{(k)}}{\sqrt{d}} + \left(\frac{\boldsymbol{Q}_{T_{j+k},:}^{C}\boldsymbol{V}_{T_{j+k},:}^{U}}{\sqrt{d}} \odot \widetilde{\boldsymbol{M}}_{T_{j+k},T_{j+k}}^{C}\right) \left(\operatorname{sigmoid}\left(\frac{\boldsymbol{Q}_{T_{j},:}^{U}\boldsymbol{K}_{T_{j+k},:}^{U}}{\sqrt{d}} + \boldsymbol{M}_{T_{j},T_{j+k}}^{U}\right)\right)^{\top}$$

Parallel scheme of forward pass



Efficient Inference with UQ-KV Cache

• Updating step: store U, Q

$$m{U}^t = \left(m{U}^{t-1} + \operatorname{sigmoid}_{m{0}^{1 imes d}} \left(m{rac{m{Q}_{t-1}^U m{k}_t^{U^{ op}}}{\sqrt{d}}}
ight) m{v}_t^U
ight)$$

Combining step: store K, V

$$oldsymbol{s}_t^C = rac{oldsymbol{q}_t^C oldsymbol{K}_t^{C^ op}}{\sqrt{d}} \in \mathbb{R}^{1 imes t}$$

$$oldsymbol{s}_t^U = rac{oldsymbol{q}_t^C oldsymbol{U}^{t^{-1}}}{\sqrt{d}} \in \mathbb{R}^{1 imes t}$$

$$\boldsymbol{p}_t = \operatorname{softmax}\left(\boldsymbol{s}_t^C - \operatorname{SiLU}\left(\boldsymbol{s}_t^U\right)\right) \in \mathbb{R}^{1 \times t}$$

attention
$$(\boldsymbol{X}^t) = \boldsymbol{p}_t \boldsymbol{V}_t^C \in \mathbb{R}^{1 \times d}$$
.

Experiment Setup

Model Architecture

$$\mathbf{Y}^{(l)} = \text{MultiHead-Attention}(\text{RMSNorm}(\mathbf{X}^{(l)}))$$

 $\mathbf{X}^{(l+1)} = \text{SwiGLU}(\text{RMSNorm}(\mathbf{Y}^{(l)})),$

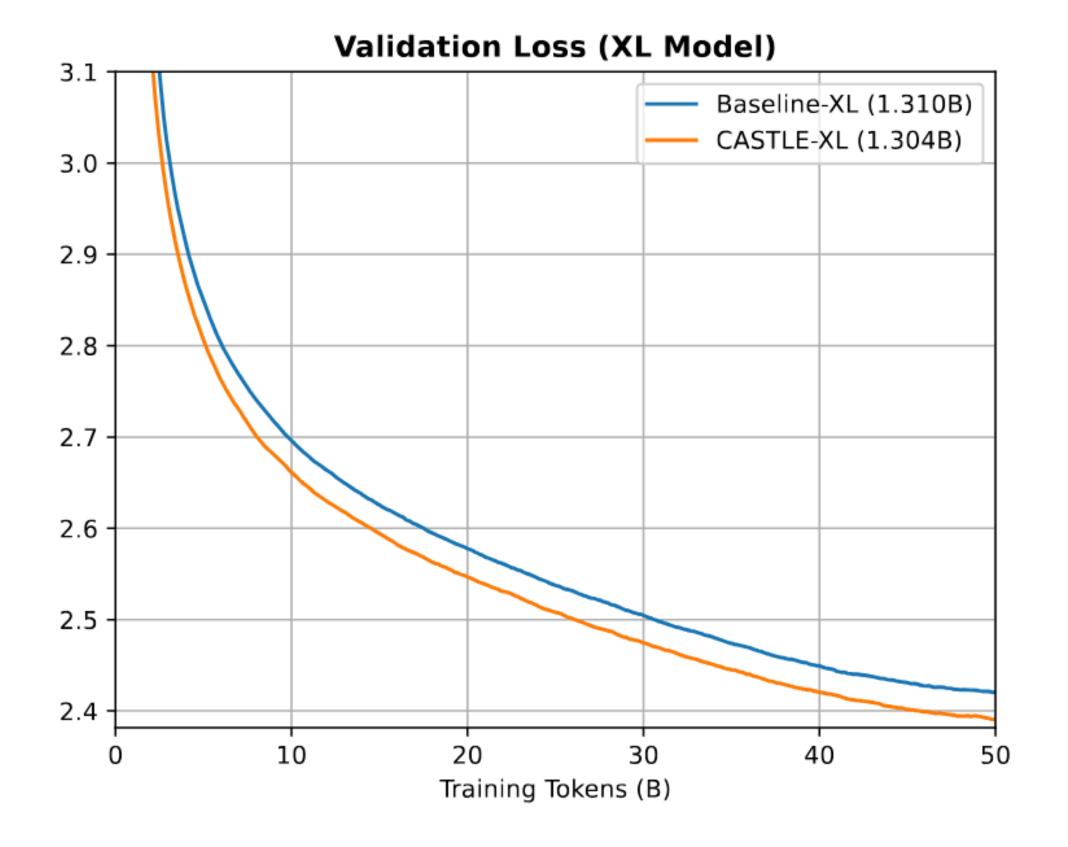
Model configuration and training recipe

Model Name	$n_{ m params}$	$n_{ m layers}$	$d_{ m model}$	$n_{ m heads}$	d	Batch Size	Learning Rate	
Baseline-S CASTLE-S	$160 \mathrm{M}$ $160 \mathrm{M}$	$\frac{12}{12}$	896 (=14 * 64) 896	14 8	64	0.5M	6.0×10^{-4}	
Baseline-M	353M	24	$1024 \ (=16 * 64)$	16	64	0.5M	3.0×10^{-4}	
CASTLE-M	351M	24	1024	9	64			
Baseline-L	756M	24	$1536 \ (=16 * 96)$	16	96	0.5M	2.5×10^{-4}	
CASTLE-L	753M	24	1536	9	96			
${\bf Baseline\text{-}XL}$	1.310B	24	$2048 \ (=16 * 128)$	16	128	0.5M	2.0×10^{-4}	
CASTLE-XL	1.304B	24	2048	9	128	0.0101	2.0 X 10	

Experimental Results

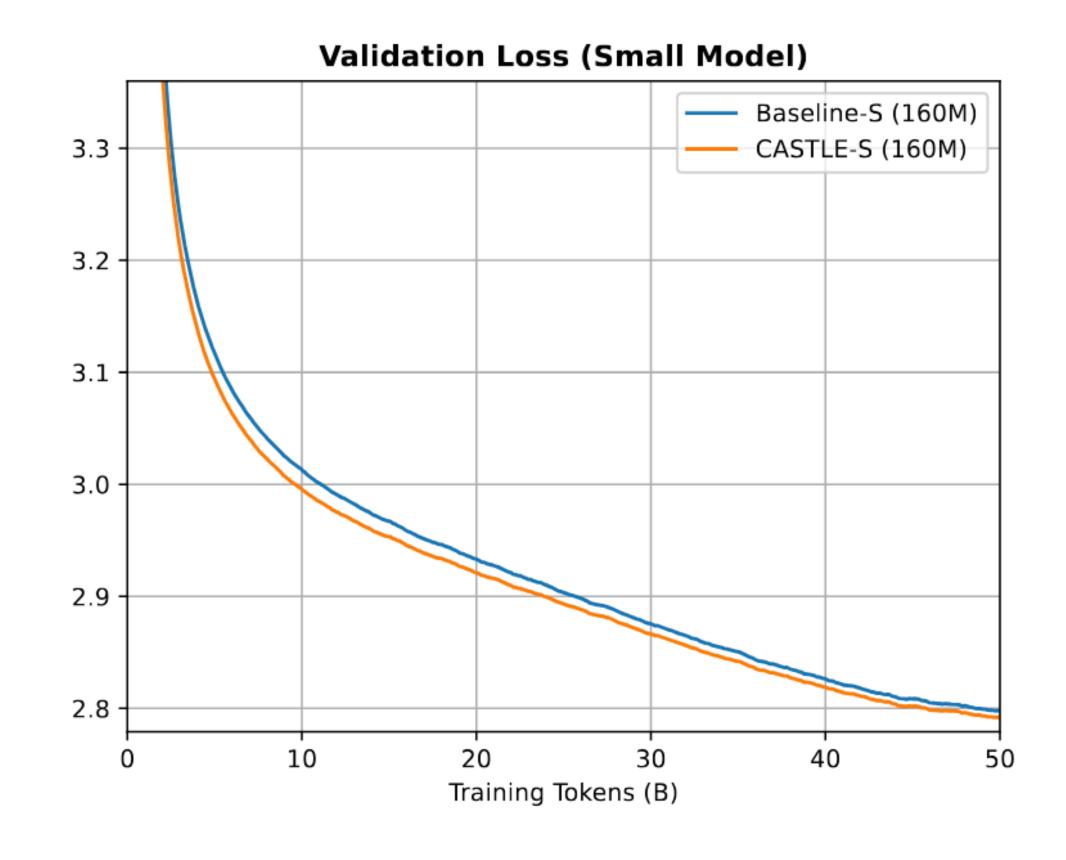
Training & Validation Loss

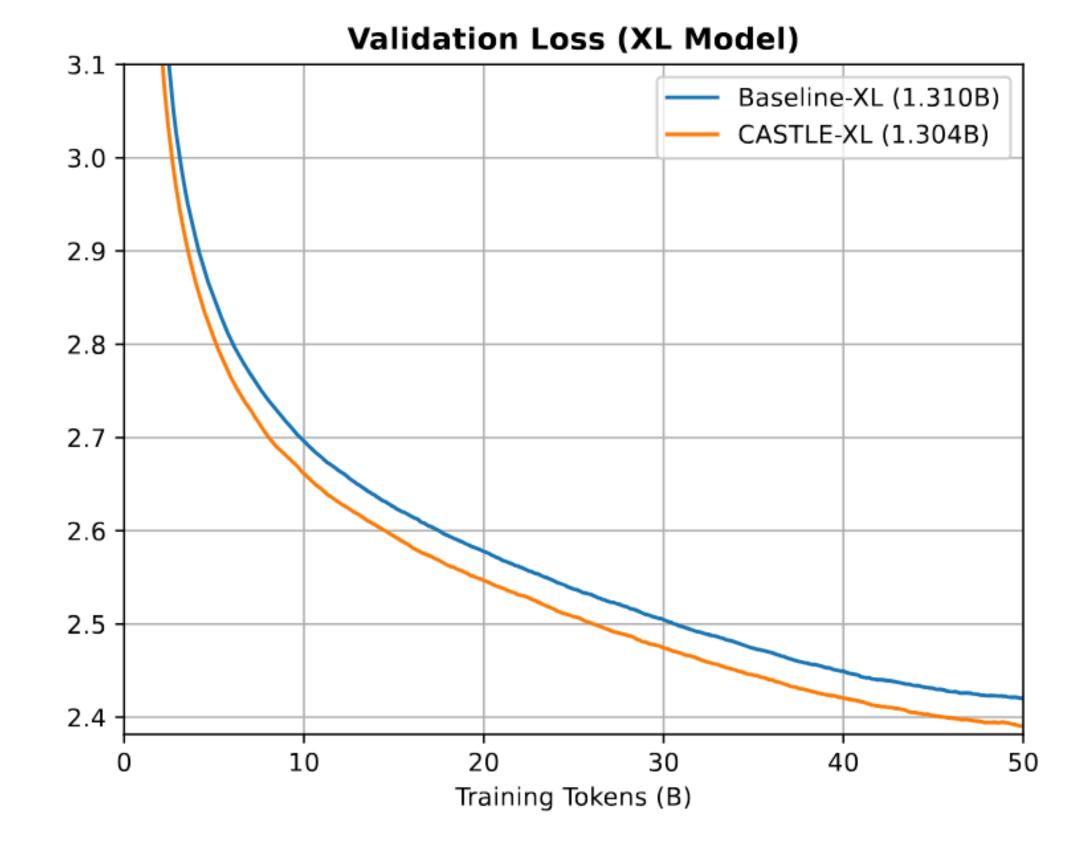
	Tr	rain	Eval			
	Loss	PPL	Loss	PPL		
Baseline-S CASTLE-S	$2.795 \\ 2.789$	16.364 16.259	$2.798 \\ 2.792$	16.411 16.315		
Baseline-M	2.641	14.030	2.639	14.004		
CASTLE-M	2.616	13.684	2.615	13.665		
Baseline-L	2.513	12.346	2.507	12.269		
CASTLE-L	2.476	11.897	2.472	11.840		
Baseline-XL	2.430	11.360	2.426	11.309		
CASTLE-XL	2.401	11.031	2.391	10.922		



Experimental Results

Comparison between loss curves of small models and XL models





Experimental Results: CASTLE vs. CASTLE-SWL

• CASTLE-SWL matches CASLTE on model scales (XL model as example here)

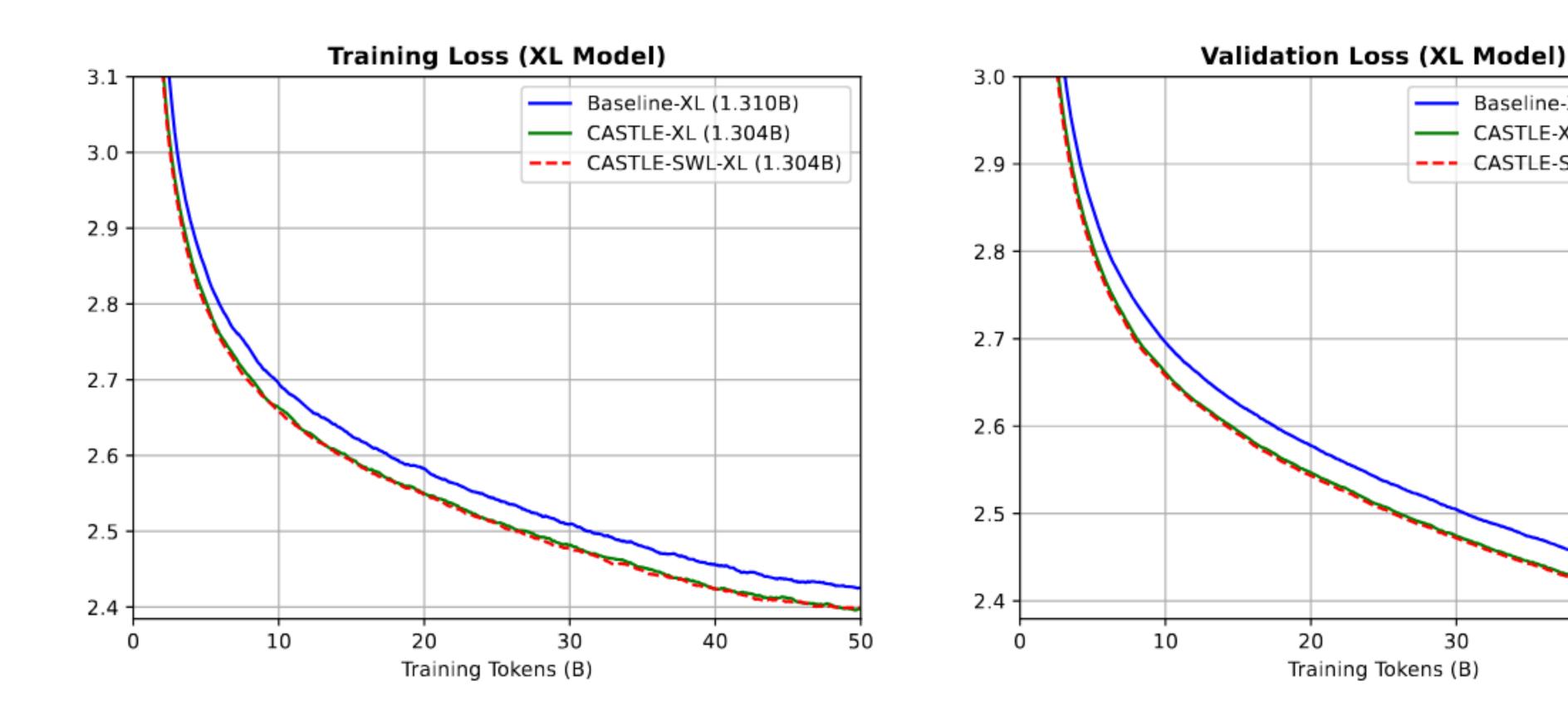
Baseline-XL (1.310B)

CASTLE-XL (1.304B)

CASTLE-SWL-XL (1.304B)

40

50



Experimental Results

Downstream Tasks (0-shot)

Model Name	ARC-C	ARC-E	BoolQ	Hella.	MMLU	OBQA	PIQA	Wino.	Avg.
Baseline-S	26.71	54.76	52.51	35.78	22.89	30.40	63.98	52.57	42.45
CASTLE-S	26.19	56.69	59.85	36.28	23.00	31.60	64.25	52.25	43.76
CASTLE-SWL-S	26.62	54.12	59.60	35.97	22.87	32.60	64.31	50.51	43.33
Baseline-M	28.58	60.90	53.61	43.01	23.21	33.40	67.95	50.91	45.20
CASTLE-M	30.20	61.36	58.01	43.24	25.34	34.60	67.95	52.64	46.67
CASTLE-SWL-M	29.52	61.41	57.03	43.76	23.29	33.00	$\boldsymbol{68.12}$	53.83	46.24
Baseline-L	32.59	65.07	57.49	47.45	23.57	32.60	70.51	50.75	47.50
CASTLE-L	32.34	$\boldsymbol{65.15}$	57.65	47.87	24.51	35.60	70.78	53.51	48.43
CASTLE-SWL-L	32.76	63.51	60.95	48.50	23.72	36.00	70.02	54.85	48.79
Baseline-XL	33.79	66.62	61.04	51.40	26.72	36.20	72.58	54.06	50.30
CASTLE-XL	35.32	67.51	62.81	52.15	23.74	37.00	70.67	56.59	50.72
CASTLE-SWL-XL	36.43	$\boldsymbol{69.07}$	60.24	51.99	24.47	37.40	71.27	55.09	50.74

Experimental Results

Downstream tasks (5-shot)

Model Name	ARC-C	ARC-E	BoolQ	Hella.	MMLU	OBQA	PIQA	Wino.	Avg.
Baseline-S	25.68	54.97	56.09	33.81	25.54	28.20	63.98	52.57	42.60
CASTLE-S	26.02	54.25	57.13	35.24	25.22	29.80	64.53	50.99	42.90
CASTLE-SWL-S	27.39	56.19	56.09	35.46	25.34	30.20	64.85	51.22	43.34
Baseline-M	31.06	62.46	48.47	42.83	25.22	33.00	68.39	51.78	45.40
CASTLE-M	32.17	64.06	$\bf 54.62$	43.47	25.22	33.80	69.48	52.49	46.91
CASTLE-SWL-M	32.85	63.85	54.19	44.18	26.43	33.80	69.04	53.43	47.22
Baseline-L	33.36	63.64	59.24	46.16	26.82	33.40	69.53	54.06	48.28
CASTLE-L	37.37	67.89	50.95	47.71	26.11	34.20	70.18	54.06	48.56
CASTLE-SWL-L	36.26	65.53	$\boldsymbol{60.58}$	48.55	24.70	34.80	69.10	53.67	49.15
Baseline-XL	35.58	65.78	61.07	50.84	26.71	36.20	71.27	52.72	50.02
CASTLE-XL	39.08	70.24	$\bf 62.60$	51.63	24.16	37.40	71.00	58.33	51.80
CASTLE-SWL-XL	38.99	70.08	61.74	52.35	25.85	37.20	$\boldsymbol{72.52}$	56.75	51.93

Ablation Studies on Number of Keys

- Is the improvement from more keys?
 - Model configurations

Model Name	$n_{ m params}$	$n_{ m layers}$	$d_{ m model}$	$n_{ m heads}$	$n_{ m Lookahead Keys} + n_{ m Causal Keys}$	d
Baseline-M	353M	24	$1024 \ (=16 * 64)$	16	16	64
CASTLE-M	351M	24	1024	9	18	64
CASTLE-M-16	340M	24	1024	8	16	64
Baseline-XL	1.310B	24	2048 (=16 * 128)	16	16	128
CASTLE-XL	1.304B	24	2048	9	18	128
CASTLE-XL-16	1.260B	24	2048	8	16	128

Training & Validation loss

		Train		Eval		Eval				T	rain	E	val
	$n_{ m params}$	Loss	PPL	Loss	PPL		$n_{ m params}$	Loss	PPL	Loss	PPL		
Baseline-M CASTLE-M CASTLE-M-16	353M 351M 340M	2.740 2.709 2.714	15.483 15.018 15.093	2.742 2.711 2.716	15.523 15.039 15.126	Baseline-M CASTLE-SWL-M CASTLE-SWL-M-16	353M 351M 340M	2.740 2.710 2.716	15.483 15.036 15.117	2.742 2.713 2.718	15.523 15.068 15.150		
Baseline-XL CASTLE-XL CASTLE-XL-16	1.310B 1.304B 1.260B	2.548 2.507 2.514	12.779 12.267 12.349	2.543 2.503 2.511	12.723 12.219 12.316	Baseline-XL CASTLE-SWL-XL CASTLE-SWL-XL-16	1.310B 1.304B 1.260B	2.548 2.506 2.513	12.779 12.255 12.339	2.543 2.503 2.508	$12.723 \\ 12.217 \\ \underline{12.276}$		

Thank you!