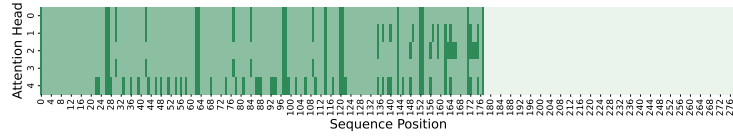


429 **Appendix**

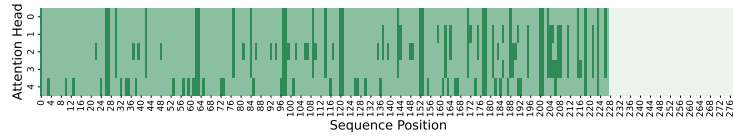
430 **A More Observation Plots**

431 **A.1 Repetitive Attention Pattern**

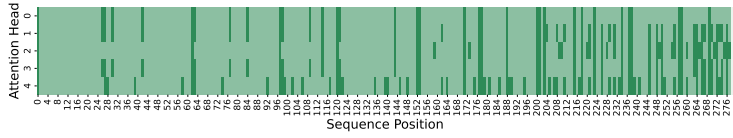
432 We provide the attention map similar to Figure 1 but from a different transformer layer on the same
433 text in Figure 6, Figure 7, Figure 8 and Figure 9. A repetitive pattern and attention sparsity can be
434 observed across layers.



(a) Attention map at position 178

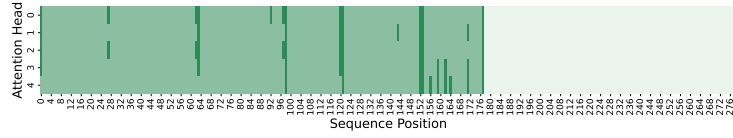


(b) Attention map at position 228

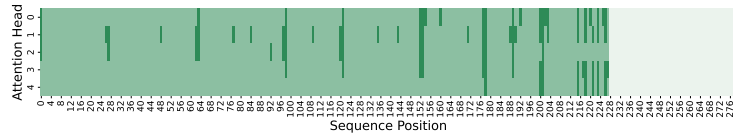


(c) Attention map at position 278

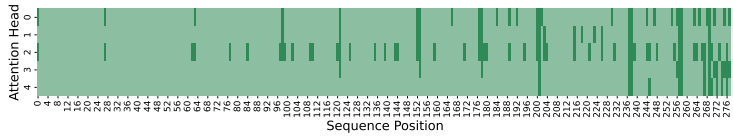
Figure 6: Attention Map at Layer 5



(a) Attention map at position 178

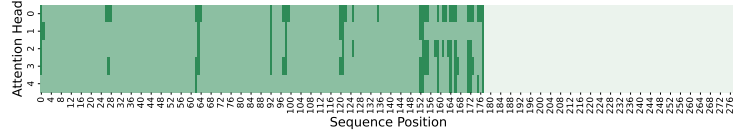


(b) Attention map at position 228

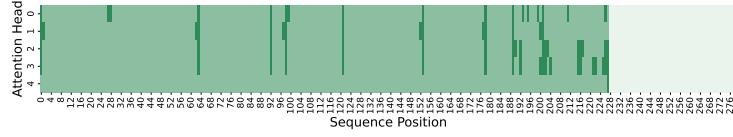


(c) Attention map at position 278

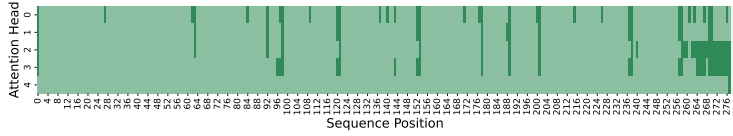
Figure 7: Attention Map at Layer 10



(a) Attention map at position 178

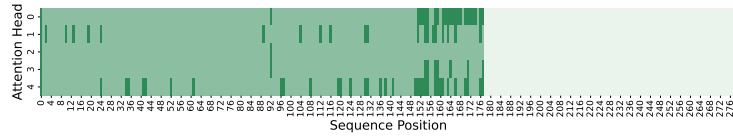


(b) Attention map at position 228

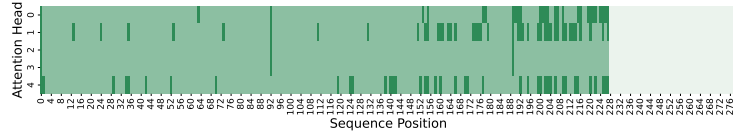


(c) Attention map at position 278

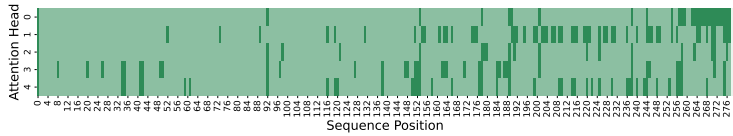
Figure 8: Attention Map at Layer 15



(a) Attention map at position 178



(b) Attention map at position 228

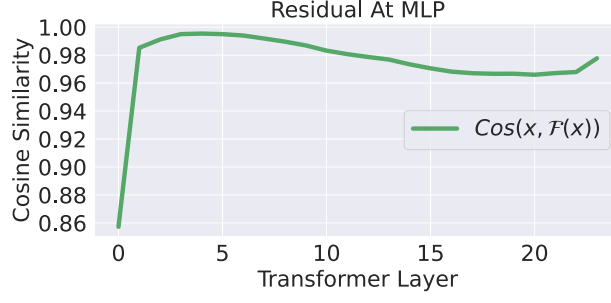


(c) Attention map at position 278

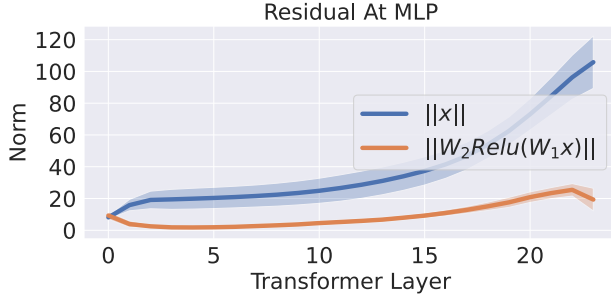
Figure 9: Attention Map at Layer 20

435 **A.2 Cross Layer Cosine Similarity**

436 In Section 3.3, our analysis assumes a large cosine similarity between the input and output of \mathcal{F} . Here,
 437 we provide empirical evidence to support such an assumption in Figure 10. Because of the residual
 438 connection in \mathcal{F} and the domination of x , the cosine similarity between x and $\mathcal{F}(x)$ is extremely
 439 high.



(a) Cosine Similarity



(b) Norm Comparison

Figure 10: x and $\mathcal{F}(x)$ is high in cosine similarity

440 B Proofs

441 B.1 Proof of Theorem 3.1

442 We consider the token generation process of a simplified model: a single-layer transformer model
443 with single-head attention.

$$x_{t+1} = \mathcal{F}(a_t), \text{ where } a_t = \text{softmax}(1/t \cdot x_t W_Q W_K^\top X_{t-1}^\top) X_{t-1} W_V W_O \quad (5)$$

444 $x_t \in \mathbb{R}^{1 \times d}$ is a row vector. $X_{t-1} \in \mathbb{R}^{(t-1) \times d}$ denotes the aggregation of x_1, \dots, x_{t-1} , where
445 the j th row is x_j . $W_Q, W_K, W_V \in \mathbb{R}^{d \times p}$ and $W_O \in \mathbb{R}^{p \times d}$ are the attention weights. Lastly,
446 $\mathcal{F}: \mathbb{R}^{1 \times d} \rightarrow \mathbb{R}^{1 \times d}$ denotes the MLP block following attention block, a two-layer MLP with skip
447 connections, given by

$$\mathcal{F}(x) = x + W_2 \text{relu}(W_1 x) \quad (6)$$

448 We are interested in the attention scores $\alpha_t = \text{softmax}(1/t \cdot x_t W_Q W_K^\top X_{t-1}^\top)$. Notice that $\alpha_{t,j}$ scales
449 with $x_t W_Q W_K^\top x_j^\top$. We first re-state the Theorem 3.1 below.

450 **Theorem B.1.** Let $A = W_V W_O W_Q W_K^\top$ and let $\lambda_K, \lambda_Q, \lambda_V, \lambda_O$ denote the largest singular values
451 of W_K, W_Q, W_V, W_O , respectively. Consider the transformer in (5) with normalized inputs $\|x_t\|_2 =$

452 1 for all t . Let $c, \epsilon > 0$ be constants. Assume that $a_t x_{t+1}^\top \geq (1 - \delta) \|a_t\|_2$ with $\delta \leq \left(\frac{c\epsilon}{\lambda_Q \lambda_K \lambda_V \lambda_O}\right)^2$.

453 Then for all x_ℓ satisfying $x_\ell A x_\ell^\top \geq c$ and $x_\ell A x_\ell \geq \epsilon^{-1} \max_{j \in [t], j \neq \ell} x_j A x_\ell^\top$, it holds that

$$\frac{x_\ell A x_\ell^\top}{\|a_t\|_2} (\alpha_{t,\ell} - 3\epsilon) \leq x_{t+1} W_Q W_K^\top x_j^\top \leq \frac{x_\ell A x_\ell^\top}{\|a_t\|_2} (\alpha_{t,\ell} + 3\epsilon) \quad (7)$$

454 As a preparation of the proof, we first show two lemmas.

455 **Lemma B.1.** Let $x_1, x_2 \in \mathbb{R}^{1 \times m}$ satisfies $\|x_1\|_2 = \|x_2\|_2 = 1$ and $x_1 x_2^\top \geq 1 - \delta$ for some
456 $\delta \in (0, 1)$. Then for all $y \in \mathbb{R}^{1 \times m}$ we have

$$|x_1 y^\top - x_2 y^\top| \leq \sqrt{2\delta} \|y\|_2$$

457 *Proof.* Let $x_2 = x_2^{\parallel} + x_2^{\perp}$ where

$$x_2^{\parallel} = x_1 x_2^{\top} \cdot x_1; \quad x_2^{\perp} = x_2 - x_2^{\parallel}$$

458 Then it is easy to see that $x_2^{\perp} x_1^{\top} = 0$. By the Pythagorean Theorem, we have

$$\|x_2^{\perp}\|_2^2 = \|x_2\|_2^2 - \|x_2^{\parallel}\|_2^2 = \delta(2 - \delta)$$

459 Therefore, we have

$$\begin{aligned} \|x_1 - x_2\|_2^2 &= \|(x_1 - x_2^{\parallel}) - x_2^{\perp}\|_2^2 \\ &= \|(1 - x_1 x_2^{\top}) x_1 - x_2^{\perp}\|_2^2 \\ &= (1 - x_1 x_2^{\top})^2 + \|x_2^{\perp}\|_2^2 \\ &= 2\delta \end{aligned}$$

460 Thus, the Cauchy-Schwarz inequality implies

$$|x_1 y^{\top} - x_2 y^{\top}| \leq \|x_1 - x_2\|_2 \cdot \|y\|_2 = \sqrt{2\delta} \|y\|_2$$

461

□

462 **Lemma B.2.** Let $\ell \in [t]$ be given. Suppose that $x_{\ell} A x_{\ell}^{\top} > \epsilon^{-1} |x_j A x_{\ell}^{\top}|$ for all $j \neq \ell$. Then we have

$$(\mathcal{S}(t)_{\ell} - \epsilon) x_{\ell}^{\top} a x_{\ell} \leq x_{\ell}^{\top} W_K^{\top} W_Q a_t \leq (\mathcal{S}(t)_{\ell} + \epsilon) x_{\ell}^{\top} a x_{\ell}$$

463 *Proof.* Notice that

$$a_t = \alpha_t X_{t-1} W_V W_O = \left(\sum_{j=1}^{t-1} \alpha_{t,j} x_j \right) W_V W_O$$

464 Thus, we have

$$a_t W_Q W_K^{\top} x_{\ell}^{\top} = \left(\sum_{j=1}^{t-1} \alpha_{t,j} x_j \right) W_V W_O W_Q W_K^{\top} x_{\ell}^{\top} = \sum_{j=1}^{t-1} \alpha_{t,j} x_j A x_{\ell}^{\top}$$

465 Therefore

$$\begin{aligned} |a_t W_Q W_K^{\top} x_{\ell}^{\top} - \alpha_{t,\ell} x_{\ell} A x_{\ell}^{\top}| &= \left| \sum_{j=1, j \neq \ell}^{t-1} \alpha_{t,j} x_j A x_{\ell}^{\top} \right| \\ &\leq \sum_{j=1, j \neq \ell}^{t-1} \alpha_{t,j} |x_j A x_{\ell}^{\top}| \\ &\leq \epsilon x_{\ell} A x_{\ell}^{\top} \sum_{j=1, j \neq \ell}^{t-1} \alpha_{t,j} \\ &\leq \epsilon x_{\ell} A x_{\ell}^{\top} \end{aligned}$$

466 where in the second inequality we use $\epsilon^{-1} |x_j A x_{\ell}^{\top}| \leq x_{\ell} A x_{\ell}^{\top}$ and in the third inequality we use

467 $\sum_{j=1, j \neq \ell}^{t-1} \alpha_{t,j} \leq \sum_{j=1}^{t-1} \alpha_{t,j} = 1$. This implies that

$$(\alpha_{t,\ell} - \epsilon) x_{\ell} A x_{\ell}^{\top} \leq a_t W_Q W_K^{\top} x_{\ell}^{\top} \leq (\alpha_{t,\ell} + \epsilon) x_{\ell} A x_{\ell}^{\top}$$

468

□

469 Now we proceed to the main body of the proof. Assume that $\|x_\ell\|_2 = 1$ for all ℓ . Using Lemma
 470 (B.1), if $a_t x_{t+1}^\top \geq (1 - \delta) \|a_t\|_2$, then we have

$$\left| \|a_t\|_2^{-1} a_t W_Q W_K^\top x_\ell^\top - x_{t+1} W_Q W_K^\top x_\ell^\top \right| \leq \sqrt{2\delta} \|W_Q W_K^\top x_\ell^\top\|_2$$

471 Recall that λ_Q, λ_K are the maximum singular values of W_Q and W_K , respectively. Then it holds
 472 that $\|W_Q W_K^\top x_\ell^\top\|_2 \leq \lambda_Q \lambda_K \|x_\ell\|_2$. Using $\|x_\ell\|_2 = 1$, we have

$$\left| \|a_t\|_2^{-1} a_t W_Q W_K^\top x_\ell^\top - x_{t+1} W_Q W_K^\top x_\ell^\top \right| \leq \sqrt{2\delta} \lambda_Q \lambda_K$$

473 Notice that

$$\begin{aligned} \|a_t\|_2 &= \left\| \left(\sum_{j=1}^{t-1} \alpha_{t,j} x_j \right) W_V W_O \right\| \\ &\leq \lambda_O \lambda_V \left\| \sum_{j=1}^{t-1} \alpha_{t,j} x_j \right\|_2 \\ &\leq \lambda_O \lambda_V \sum_{j=1}^{t-1} \alpha_{t,j} \|x_j\|_2 \\ &= \lambda_O \lambda_V \end{aligned}$$

474 Then since $\delta \leq \left(\frac{c\epsilon}{\lambda_Q \lambda_K \lambda_V \lambda_O} \right)^2$, we have

$$\left| \|a_t\|_2^{-1} a_t W_Q W_K^\top x_\ell^\top - x_{t+1} W_Q W_K^\top x_\ell^\top \right| \leq \frac{2c\epsilon}{\lambda_V \lambda_O} \leq \frac{2c\epsilon}{\|a_t\|_2}$$

475 Since by Lemma (B.2), we have

$$\left| a_t W_Q W_K^\top x_\ell^\top - \alpha_{t,\ell} x_\ell A x_\ell^\top \right| \leq \epsilon x_\ell^\top a x_\ell$$

476 It must hold that

$$\left| x_{t+1} W_Q W_K^\top x_\ell^\top - \|a_{t+1}\|_2^{-1} \alpha_{t,\ell} x_\ell A x_\ell^\top \right| \leq \frac{\epsilon}{\|a_t\|_2} x_\ell^\top a x_\ell + \frac{2c\epsilon}{\|a_t\|_2}$$

477 Since $x_\ell^\top a x_\ell \geq c$, it holds that

$$\frac{2c\epsilon}{\|a_t\|_2} \leq \frac{2\epsilon}{\|a_t\|_2} x_\ell^\top a x_\ell$$

478 which implies that

$$\left| x_{t+1} W_Q W_K^\top x_\ell^\top - \|a_t\|_2^{-1} \alpha_{t,\ell} x_\ell A x_\ell^\top \right| \leq \frac{3\epsilon}{\|a_t\|_2} x_\ell^\top a x_\ell$$

479 Therefore

$$\frac{x_\ell A x_\ell^\top}{\|a_t\|_2} (\alpha_{t,\ell} - 3\epsilon) \leq x_{t+1} W_Q W_K^\top x_\ell^\top \leq \frac{x_\ell A x_\ell^\top}{\|a_t\|_2} (\alpha_{t,\ell} + 3\epsilon)$$

480 B.2 Proof of Theorem 4.1

481 Let $\{\tilde{x}_t\}_{t=0}^T$ denote the tokens generated by the transformer with budget KV cache as in Algorithm 2
 482 with $m = 1$:

$$\tilde{x}_{t+1} = \mathcal{F}(\tilde{a}_t), \text{ where } \tilde{a}_t = \text{softmax} \left(\frac{1}{t} \cdot \tilde{x}_t W_Q \tilde{\mathcal{K}}_t^\top \right) \tilde{\mathcal{V}}_t^\top W_O$$

483 Notice that when $m = 1$, i.e., in each iteration, we drop one token with the lowest score, the cache
 484 will always maintain B tokens. If the ranking of the attention scores does not change in each iteration,
 485 Algorithm 2 will always drop tokens with the smallest attention scores.

486 For reference purposes, let $\{x_t\}_{t=0}^T$ denote the tokens generated by a vanilla transformer defined in
 487 (5). We re-state Theorem 4.1 below, which bounds the difference $\|x_t - \tilde{x}_t\|_2$.

488 **Theorem B.2.** Let λ_1, λ_2 denote the largest singular values of W_1 and W_2 in (6). Let

$$\beta_{t,j} = \frac{\exp(1/t \cdot \tilde{x}_t W_Q W_K^\top \tilde{x}_j^\top)}{\sum_{i=1}^{t-1} \exp(1/t \cdot \tilde{x}_t W_Q W_K^\top \tilde{x}_i^\top)}$$

489 and assume that each $\beta_{t,j} = cv_{t,j}$, where $v_{t,j}$ are sampled from a power-law distribution with pdf
 490 $f(x) = c(x+b)^{-k}$. Suppose that $\lambda_V \lambda_O (1 + \lambda_1 \lambda_2)(1 + \lambda_Q \lambda_K) \leq \frac{1}{2}$. Let T_{\min} and T_{\max} denote
 491 the starting and maximum sequence lengths, respectively, and let $B \leq T_{\max}$ denote the budget as
 492 in Algorithm 2. If for all $t \in [T_{\min}, T_{\max}]$, S_t contains only tokens with at most the largest B values
 493 of $\beta_{t,j}$, that is, $|S_t| = B$ and $\min_{j \in S_t} \beta_{t,j} \geq \max_{j \notin S_t} \beta_{t,j}$, then for all $\epsilon \in (0, 1)$, with probability
 494 at least $1 - T_{\max} \exp\left(-\frac{\epsilon^2 b^2 (T_{\min} - 1)}{(k-2)^2 (u-b)^2}\right) - T_{\max} \exp\left(-\frac{2(T_{\min} - 1)(1-B/T_{\max})^2}{(1-\epsilon)^2}\right)$, the following error
 495 bound must hold for all $t \in [T_{\min}, T_{\max}]$

$$\mathbb{E} [\|x_t - \tilde{x}_t\|_2] \leq \frac{2.1(1 - B/T_{\max})}{(1 - \epsilon)^2} \left(k - (k - 1) \left(\frac{1 - \epsilon}{B/T_{\max} - \epsilon} \right)^{1/(k-1)} \right)$$

496 Define $m_{k,j} = \mathbb{I}\{j \in S_t\}$. With the definition of $m_{k,j}$, \tilde{a}_t can be written as

$$\tilde{a}_t = \left(\sum_{j=1}^{t-1} \tilde{\alpha}_{t,j} \tilde{x}_j \right) W_V W_O; \quad \tilde{\alpha}_{t,j} = \frac{m_{k,j} \exp(1/t \cdot \tilde{x}_t W_Q W_K^\top \tilde{x}_j^\top)}{\sum_{i=1}^{t-1} m_{k,i} \exp(1/t \cdot \tilde{x}_t W_Q W_K^\top \tilde{x}_i^\top)} \quad (8)$$

497 Our first lemma shows the Lipschitzness of the attention module.

498 **Lemma B.3.** Consider two sequences of tokens $\{x_i\}_{i=1}^t$ and $\{y_i\}_{i=1}^t$ where $\|x_i\|_2 = \|y_i\|_2 = 1$ for
 499 all $i \in [t]$. Define $X_{t-1}, Y_{t-1} \in \mathbb{R}^{(t-1) \times d}$ as the matrices whose i th row are x_i and y_i , respectively.
 500 Let $\Delta_t = \|x_t - y_t\|_2$. Then we have

$$\left\| \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top X_{t-1}^\top \right)_2 - \text{softmax} \left(\frac{1}{t} y_t W_Q W_K^\top Y_{t-1}^\top \right)_2 \right\|_2 \leq 2 \frac{\sqrt{t-1}}{t} \lambda_Q \lambda_K \Delta_t$$

501 *Proof.* We can decompose the difference as

$$\begin{aligned} & \left\| \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top X_{t-1}^\top \right) - \text{softmax} \left(\frac{1}{t} y_t W_Q W_K^\top Y_{t-1}^\top \right) \right\|_2 \\ & \leq \left\| \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top X_{t-1}^\top \right) - \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top Y_{t-1}^\top \right) \right\|_2 \\ & \quad + \left\| \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top Y_{t-1}^\top \right) - \text{softmax} \left(\frac{1}{t} y_t W_Q W_K^\top Y_{t-1}^\top \right) \right\|_2 \end{aligned}$$

502 By the Lipschitzness of softmax, we have

$$\begin{aligned} & \left\| \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top X_{t-1}^\top \right) - \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top Y_{t-1}^\top \right) \right\|_2 \\ & \leq \frac{1}{t} \left\| x_t W_Q W_K^\top (X_{t-1} - Y_{t-1})^\top \right\|_2 \\ & \leq \frac{1}{t} \lambda_Q \lambda_K \|x_t\|_2 \|X_{t-1} - Y_{t-1}\|_2 \end{aligned}$$

503 Since $\|x_t\|_2 = 1$ and $\|X_{t-1} - Y_{t-1}\|_2 = \left(\sum_{j=1}^{t-1} \|x_j - y_j\|_2^2 \right)^{\frac{1}{2}} \leq \sqrt{t-1} \Delta_t$, we have

$$\left\| \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top X_{t-1}^\top \right) - \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top Y_{t-1}^\top \right) \right\|_2 \leq \frac{\sqrt{t-1}}{t} \lambda_Q \lambda_K \Delta_t$$

504 Similarly,

$$\begin{aligned} & \left\| \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top Y_{t-1}^\top \right) - \text{softmax} \left(\frac{1}{t} y_t W_Q W_K^\top Y_{t-1}^\top \right) \right\|_2 \\ & \leq \frac{1}{t} \left\| (x_t - y_t) W_Q W_K^\top Y_{t-1}^\top \right\|_2 \\ & \leq \frac{1}{t} \lambda_Q \lambda_K \|Y_{t-1}\|_F \|x_t - y_t\|_2 \end{aligned}$$

505 Since $\|x_t - y_t\|_2 = \Delta_t$ and $\|Y_{t-1}\|_2 = \sqrt{t-1}$, we have

$$\left\| \text{softmax} \left(\frac{1}{t} x_t W_Q W_K^\top Y_{t-1}^\top \right) - \text{softmax} \left(\frac{1}{t} y_t W_Q W_K^\top Y_{t-1}^\top \right) \right\|_2 \leq \frac{\sqrt{t-1}}{t} \lambda_Q \lambda_K \Delta_t$$

506 Combining the two bounds gives

$$\left\| \text{softmax} \left(\frac{1}{\sqrt{t}} x_t W_Q W_K^\top X_{t-1}^\top \right) - \text{softmax} \left(\frac{1}{\sqrt{t}} y_t W_Q W_K^\top Y_{t-1}^\top \right) \right\|_2 \leq 2 \frac{\sqrt{t-1}}{t} \lambda_Q \lambda_K \Delta_t$$

507 □

508 Our second lemma shows the difference between the output of the sampled and vanilla transformer
509 when the input is the same.

510 **Lemma B.4.** Let \tilde{a}_t be defined as in [\(8\)](#). Define b_t as

$$b_t = \left(\sum_{j=1}^{t-1} \beta_{t,j} \tilde{x}_j \right) W_V W_O; \quad \beta_{t,j} = \frac{\exp(1/t \cdot \tilde{x}_t W_Q W_K^\top \tilde{x}_j^\top)}{\sum_{i=1}^{t-1} \exp(1/t \cdot \tilde{x}_t W_Q W_K^\top \tilde{x}_i^\top)} \quad (9)$$

511 Assume that $\|x_j\|_2 = 1$ for all $j \in [t]$. Then we have

$$\|\tilde{a}_t - b_t\|_2 \leq \lambda_V \lambda_O \sum_{j \notin \hat{S}_t} \beta_{t,j}$$

512 *Proof.* A direction computation yields

$$\tilde{a}_t - b_t = \left(\sum_{j=1}^{t-1} (\tilde{\alpha}_{t,j} - \beta_{t,j}) \tilde{x}_j \right) W_V W_O$$

513 Thus, $\|\tilde{a}_t - b_t\|_2$ can be bounded as

$$\|\tilde{a}_t - b_t\|_2 \leq \lambda_V \lambda_O \sum_{j=1}^{t-1} (\tilde{\alpha}_{t,j} - \beta_{t,j}) \|\tilde{x}_j\|_2 = \lambda_V \lambda_O \sum_{j=1}^{t-1} (\tilde{\alpha}_{t,j} - \beta_{t,j})$$

514 since $\|\tilde{x}_j\|_2 = 1$ for all $j \in [t]$. Now we analyze $\tilde{\alpha}_{t,j} - \beta_{t,j}$. Let $\hat{S}_t = S_t \setminus \{t\}$. Then $m_{k,j} = 1$ if
515 and only if $j \in \hat{S}_t$. For convenience, let $r_{t,j} = 1/t \cdot \tilde{x}_t W_Q W_K^\top \tilde{x}_j^\top$. Thus, β can be written as

$$\beta_{t,j} = \frac{\exp(r_{t,j})}{\sum_{i \in \hat{S}_t} \exp(r_{t,i}) + \sum_{i \notin \hat{S}_t} \exp(r_{t,i})}$$

516 Furthermore, for all $j \notin \hat{S}_t$, we have $\tilde{\alpha}_{t,j} = 0$. For all $j \in \hat{S}_t$, we have

$$\tilde{\alpha}_{t,j} = \frac{\exp(r_{t,j})}{\sum_{i \in \hat{S}_t} \exp(r_{t,i})}$$

517 Therefore, for all $j \in \hat{S}_t$, we have

$$\begin{aligned} \beta_{t,j} - \tilde{\alpha}_{t,j} &= \exp(r_{t,j}) \cdot \frac{\sum_{i \notin \hat{S}_t} \exp(r_{t,i})}{\left(\sum_{i \in \hat{S}_t} \exp(r_{t,i}) \right) \left(\sum_{i \in \hat{S}_t} \exp(r_{t,i}) + \sum_{i \notin \hat{S}_t} \exp(r_{t,i}) \right)} \\ &= \frac{\exp(r_{t,j})}{\sum_{i \in \hat{S}_t} \exp(r_{t,i})} \cdot \frac{\sum_{i \notin \hat{S}_t} \exp(r_{t,i})}{\sum_{i \in \hat{S}_t} \exp(r_{t,i}) + \sum_{i \notin \hat{S}_t} \exp(r_{t,i})} \\ &= \tilde{\alpha}_{t,j} \sum_{i \notin \hat{S}_t} \beta_{t,j} \end{aligned}$$

518 Therefore, the bound of $\|\tilde{a}_t - b_t\|_2$ can be written as

$$\|\tilde{a}_t - b_t\|_2 \leq \lambda_V \lambda_O \left(\sum_{j \in \hat{S}_t} \tilde{\alpha}_{t,j} \sum_{i \notin \hat{S}_t} \beta_{t,j} - \sum_{j \notin \hat{S}_t} \beta_{t,j} \right) = 2\lambda_V \lambda_O \sum_{j \notin \hat{S}_t} \beta_{t,j}$$

519 where the last equality follows from $\sum_{j \in \hat{S}_t} \tilde{\alpha}_{t,j} = 1$. □

520 Our last lemma shows the Lipschitzness of the MLP in (6).

521 **Lemma B.5.** Let λ_1, λ_2 denote the largest singular values of W_1, W_2 in (6). For all $x_1, x_2 \in \mathbb{R}^d$,
 522 we have

$$\|\mathcal{F}(x_1) - \mathcal{F}(x_2)\| \leq (1 + \lambda_1 \lambda_2) \|x_1 - x_2\|_2$$

523 *Proof.* Direct computation yields

$$\begin{aligned} \|\mathcal{F}(x_1) - \mathcal{F}(x_2)\| &= \|(x_1 + W_2 \text{relu}(W_1 x_1)) - (x_2 + W_2 \text{relu}(W_1 x_2))\| \\ &\leq \|x_1 - x_2\|_2 + \|W_2 \text{relu}(W_1 x_1) - W_2 \text{relu}(W_1 x_2)\| \\ &\leq \|x_1 - x_2\|_2 + \lambda_2 \|\text{relu}(W_1 x_1) - \text{relu}(W_1 x_2)\| \\ &\leq \|x_1 - x_2\|_2 + \lambda_2 \|W_1(x_1 - x_2)\|_2 \\ &\leq \|x_1 - x_2\|_2 + \lambda_1 \lambda_2 \|x_1 - x_2\|_2 \\ &= (1 + \lambda_1 \lambda_2) \|x_1 - x_2\|_2 \end{aligned}$$

524 where in the third inequality we use the fact that $\text{relu}(\cdot)$ is 1-Lipschitz. □

525 Now we turn to the proof of our main theorem. Combining all of the results, we have

$$\begin{aligned} a_t - \tilde{a}_t &= \left(\sum_{j=1}^{t-1} \alpha_{t,j} x_j \right) W_V W_O - \left(\sum_{j=1}^{t-1} \tilde{\alpha}_{t,j} \tilde{x}_j \right) W_V W_O \\ &= \underbrace{\left(\sum_{j=1}^{t-1} \alpha_{t,j} x_j \right) W_V W_O - \left(\sum_{j=1}^{t-1} \alpha_{t,j} \tilde{x}_j \right) W_V W_O}_{\mathcal{T}_1} \\ &\quad + \underbrace{\left(\sum_{j=1}^{t-1} \alpha_{t,j} \tilde{x}_j \right) W_V W_O - \left(\sum_{j=1}^{t-1} \beta_{t,j} \tilde{x}_j \right) W_V W_O}_{\mathcal{T}_2} \\ &\quad + \underbrace{\left(\sum_{j=1}^{t-1} \beta_{t,j} \tilde{x}_j \right) W_V W_O - \left(\sum_{j=1}^{t-1} \tilde{\alpha}_{t,j} \tilde{x}_j \right) W_V W_O}_{\mathcal{T}_3} \end{aligned}$$

526 Therefore, by triangle inequality, we have

$$\|a_t - \tilde{a}_t\|_2 \leq \|\mathcal{T}_1\|_2 + \|\mathcal{T}_2\|_2 + \|\mathcal{T}_3\|_2 \tag{10}$$

527 To start, the magnitude of \mathcal{T}_1 can be bounded as

$$\begin{aligned} \|\mathcal{T}_1\|_2 &= \left\| \left(\sum_{j=1}^{t-1} \alpha_{t,j} (x_{t,j} - \tilde{x}_{t,j}) \right) W_V W_O \right\|_2 \\ &\leq \lambda_V \lambda_O \left\| \sum_{j=1}^{t-1} \alpha_{t,j} (x_{t,j} - \tilde{x}_{t,j}) \right\| \\ &\leq \lambda_V \lambda_O \sum_{j=1}^{t-1} \alpha_{t,j} \|x_{t,j} - \tilde{x}_{t,j}\|_2 \\ &\leq \lambda_V \lambda_O \Delta_t \sum_{j=1}^{t-1} \alpha_{t,j} \\ &= \lambda_V \lambda_O \Delta_t \end{aligned}$$

528 where in the third inequality we use $\|x_{t,j} - \tilde{x}_{t,j}\|_2 = \Delta_t$ and in the last equality we use
 529 $\sum_{j=1}^{t-1} \alpha_{t,j} = 1$. To bound the magnitude of \mathcal{T}_2 , we apply Lemma [B.3](#), which shows that
 530 $\|\alpha_t - \beta_t\| \leq 2 \frac{\sqrt{t-1}}{t} \lambda_Q \lambda_K \Delta_t$ to get that

$$\begin{aligned} \|\mathcal{T}_2\|_2 &= \left\| \left(\sum_{j=0}^{t-1} (\alpha_{t,j} - \beta_{t,j}) \tilde{x}_j \right) W_V W_O \right\|_2 \\ &\leq \lambda_V \lambda_O \left\| \left(\sum_{j=0}^{t-1} (\alpha_{t,j} - \beta_{t,j}) \tilde{x}_j \right) \right\|_2 \\ &\leq \lambda_V \lambda_O \sum_{j=0}^{t-1} |\alpha_{t,j} - \beta_{t,j}| \|\tilde{x}_j\|_2 \\ &\leq \lambda_V \lambda_O \|\alpha_t - \beta_t\|_1 \\ &\leq \sqrt{t-1} \lambda_V \lambda_O \|\alpha_t - \beta_t\|_2 \\ &\leq 2 \left(1 - \frac{1}{t} \right) \lambda_Q \lambda_K \lambda_V \lambda_O \Delta_t \end{aligned}$$

531 Lastly, to bound the magnitude of \mathcal{T}_3 , we use Lemma [B.4](#) to get that

$$\|\mathcal{T}_3\|_2 \leq 2 \lambda_V \lambda_O \sum_{j \notin \hat{S}_t} \beta_{t,j}$$

532 Putting things together for [\(10\)](#), we have

$$\|a_t - \tilde{a}_t\|_2 \leq \lambda_V \lambda_O \left(2 \sum_{j \notin \hat{S}_t} \beta_{t,j} + (2 \lambda_Q \lambda_K + 1) \Delta_t \right)$$

533 By Lemma [B.5](#) we can further show that

$$\|x_{t+1} - \tilde{x}_{t+1}\|_2 \leq (1 + \lambda_1 \lambda_2) \lambda_V \lambda_O \left(2 \sum_{j \notin \hat{S}_t} \beta_{t,j} + (2 \lambda_Q \lambda_K + 1) \Delta_t \right)$$

534 By Theorem [B.3](#), we have that with probability at least $1 - T_{\max} \exp\left(-\frac{\epsilon^2 b^2 (T_{\min}-1)}{(k-2)^2 (u-b)^2}\right) -$
 535 $T_{\max} \exp\left(-\frac{2(T_{\min}-1)(1-B/T_{\max})^2}{(1-\epsilon)^2}\right)$, it holds for all $t \in [T_{\min}, T_{\max}]$ that

$$\mathbb{E} \left[\sum_{j \notin \hat{S}_t} \beta_{t,j} \right] \leq \frac{(1 - B/T_{\max})}{0.98(1 - \epsilon)^2} \left(k - (k-1) \left(\frac{1 - \epsilon}{B/T_{\max} - \epsilon} \right)^{\frac{1}{k-1}} \right) := \Delta_{\max}$$

536 Given that $\mathbb{E}[\|x_t - \tilde{x}_t\|] \leq 2 \Delta_{\max}$, we have

$$\begin{aligned} \mathbb{E}[\|x_{t+1} - \tilde{x}_{t+1}\|_2] &\leq (1 + \lambda_1 \lambda_2) \lambda_V \lambda_O (2 \Delta_{\max} + 2 (2 \lambda_Q \lambda_K + 1) \Delta_{\max}) \\ &\leq 4 \lambda_V \lambda_O (1 + \lambda_1 \lambda_2) (1 + \lambda_Q \lambda_K) \Delta_{\max} \end{aligned}$$

537 Thus, as long as $\lambda_V \lambda_O (1 + \lambda_1 \lambda_2) (1 + \lambda_Q \lambda_K) \leq \frac{1}{2}$, we can guarantee that

$$\mathbb{E}[\|x_{t+1} - \tilde{x}_{t+1}\|_2] \leq 2 \Delta_{\max}$$

538 Thus, for all $t \in [T_{\min}, T_{\max}]$, we have that

$$\mathbb{E}[\|x_t - \tilde{x}_t\|_2] \leq \frac{2.1(1 - B/T_{\max})}{(1 - \epsilon)^2} \left(k - (k-1) \left(\frac{1 - \epsilon}{B/T_{\max} - \epsilon} \right)^{\frac{1}{k-1}} \right)$$

539 **B.3 Budgeted Cache**

540 **Theorem B.3.** Let $\beta_{t,j}$ be sampled from some power-law distribution $f(x) = c(x+b)^{-\gamma}$ with support
 541 on $[0, u-b]$ for some $k > 2$ and $u \geq 5b$. Let S_t be defined in Theorem B.2 and define $\hat{S}_t = S_t \setminus \{t\}$.
 542 Then with probability at least $1 - T_{\max} \exp\left(-\frac{\epsilon^2 b^2 (T_{\min}-1)}{(k-2)^2 (u-b)^2}\right) - T_{\max} \exp\left(-\frac{2(T_{\min}-1)(1-B)^2}{(1-\epsilon)^2}\right)$ it
 543 holds for all $t \in T$ that

$$\mathbb{E} \left[\sum_{j \notin \hat{S}_t} \beta_{t,j} \right] \leq \frac{(1 - B/T_{\max})}{0.98(1-\epsilon)^2} \left(k - (k-1) \left(\frac{1-\epsilon}{B/T_{\max} - \epsilon} \right)^{\frac{1}{k-1}} \right) \quad (11)$$

544 We consider the case of maintaining a budget of B by dropping the smallest $\beta_{t,j}$'s. Assume that v_j
 545 has pdf $f(x) = c(x+b)^{-k}$ with support on $[0, u-b]$. To make things precise, we first compute c

$$c = \left(\int_0^{u-b} (x+b)^{-k} dx \right)^{-1} = \frac{k-1}{b^{1-k} - u^{1-k}}$$

546 To start, we notice that

$$\int x(x+b)^{-k} = -\frac{(x+b)^{1-k}((k-1)x+b)}{(k-1)(k-2)} := g(x)$$

547 Let $C = \sum_{j=1}^{t-1} v_j$, then the expectation of C is

$$\begin{aligned} \mathbb{E}[C] &= (t-1)\mathbb{E}[v_1] = (t-1) \frac{k-1}{b^{1-k} - u^{1-k}} \int_0^\infty x(x+b)^{-k} dx \\ &= (t-1) \frac{k-1}{b^{1-k} - u^{1-k}} (g(u) - g(0)) \\ &= (t-1) \frac{k-1}{b^{1-k} - u^{1-k}} \left(\frac{b^{2-k}}{(k-1)(k-2)} - \frac{u^{1-k}((k-1)u - (k-2)b)}{(k-1)(k-2)} \right) \\ &= \frac{t-1}{k-2} \cdot \frac{b^{2-k} - (k-1)u^{2-k} + (k-2)bu^{1-k}}{b^{1-k} - u^{1-k}} \end{aligned}$$

548 Let $\Delta = \frac{b^{2-k} - (k-1)u^{2-k} + (k-2)bu^{1-k}}{b^{1-k} - u^{1-k}}$. By Hoeffding's inequality, we have that

$$\mathbb{P}(C \leq (1-\epsilon)\mathbb{E}[C]) \leq \exp\left(-\frac{2\epsilon^2 \mathbb{E}[C]^2}{(t-1)(u-b)^2}\right)$$

549 This implies that with probability at least $1 - \exp\left(-\frac{2\epsilon^2 \Delta^2 (t-1)}{(k-2)^2 (u-b)^2}\right)$ we have

$$C \geq (1-\epsilon)\Delta \frac{t-1}{k-2}$$

550 Now, we proceed to bound $\sum_{j \notin \hat{S}_t} \beta_{t,j}$ where $\hat{S}_t = \{j \in [t-1] : \beta_{t,j} \geq \frac{\gamma}{C}\}$. Equivalently, we can
 551 bound $C^{-1} \sum_{j=1}^{t-1} \mathbb{I}\{v_j \leq \gamma\} v_j$. Its expectation is given by

$$\begin{aligned} \mathbb{E} \left[C^{-1} \sum_{j=1}^{t-1} \mathbb{I}\{v_j \leq \gamma\} v_j \right] &\leq \frac{k-2}{(t-1)\Delta(1-\epsilon)} \mathbb{E} \left[\sum_{j=1}^{t-1} \mathbb{I}\{v_j \leq \gamma\} v_j \right] \\ &= \frac{k-2}{\Delta(1-\epsilon)} \cdot \frac{k-1}{b^{1-k} - u^{1-k}} \int_0^\gamma x(x+b)^{-k} dx \\ &= \frac{(k-1)(k-2)}{\Delta(1-\epsilon)(b^{1-k} - u^{1-k})} (g(\gamma) - g(0)) \end{aligned}$$

552 We pause here and study how small can we choose γ . Notice that

$$\mathbb{E} \left[\sum_{j=1}^{t-1} \mathbb{I}\{v_j \leq \gamma\} \right] = (t-1)\mathbb{P}(v_j \leq \gamma) = (t-1) \cdot \frac{b^{1-k} - (\gamma+b)^{1-k}}{b^{1-k} - u^{1-k}}$$

553 By Hoeffding's inequality again, we have that

$$\begin{aligned} & \mathbb{P} \left(\sum_{j=1}^{t-1} \mathbb{I}\{v_j \leq \gamma\} \geq (1-\epsilon)(t-1) \cdot \frac{b^{1-k} - (\gamma+b)^{1-k}}{b^{1-k} - u^{1-k}} \right) \\ & \leq \exp \left(- \frac{2(t-1)\epsilon^2 (b^{1-k} - (\gamma+b)^{1-k})^2}{(b^{1-k} - u^{1-k})^2} \right) \end{aligned}$$

554 Enforcing $\sum_{j=1}^{t-1} \mathbb{I}\{v_j \leq \gamma\} \geq T_{\max} - B$ gives $(\gamma+b)^{1-k} \leq b^{1-k} - \frac{1-B/T_{\max}}{1-\epsilon} (b^{1-k} - u^{1-k})$,

555 which can be satisfied as long as $\gamma \geq \left(\left(\frac{B/T_{\max} - \epsilon}{1-\epsilon} \right)^{\frac{1}{1-k}} - 1 \right) b$. Therefore

$$g(\gamma) = - \left(b^{1-k} - \frac{1-B/T_{\max}}{1-\epsilon} (b^{1-k} - u^{1-k}) \right) \frac{b + (k-1)\gamma}{(k-1)(k-2)}$$

556 We further notice that

$$b^{1-k} - \frac{1-B/T_{\max}}{1-\epsilon} (b^{1-k} - u^{1-k}) \geq \frac{B/T_{\max} - \epsilon}{1-\epsilon} (b^{1-k} - u^{1-k})$$

557 This gives

$$\begin{aligned} \mathbb{E} \left[C^{-1} \sum_{j=1}^{t-1} \mathbb{I}\{v_j \leq \gamma\} v_j \right] & \leq \frac{b(1-B/T_{\max})}{\Delta(1-\epsilon)^2} - \frac{(k-1)(B/T_{\max} - \epsilon)\gamma}{\Delta(1-\epsilon)^2} \\ & \leq \frac{b(1-B/T_{\max})}{\Delta(1-\epsilon)^2} \left(k - (k-1) \left(\frac{1-\epsilon}{B/T_{\max} - \epsilon} \right)^{\frac{1}{k-1}} \right) \end{aligned}$$

558 Notice that if $u \geq 5b$, we have

$$\Delta = b - (k-1) \left(\frac{u}{b} \right)^{1-k} \cdot \frac{b-u}{b^{1-k} - u^{1-k}} \leq 0.98b$$

559 Therefore

$$\mathbb{E} \left[C^{-1} \sum_{j=1}^{t-1} \mathbb{I}\{v_j \leq \gamma\} v_j \right] \leq \frac{(1-B/T_{\max})}{0.98(1-\epsilon)^2} \left(k - (k-1) \left(\frac{1-\epsilon}{B/T_{\max} - \epsilon} \right)^{\frac{1}{k-1}} \right)$$

560 holds with probability at least $1 - \exp \left(- \frac{\epsilon^2 b^2 (t-1)}{(k-2)^2 (u-b)^2} \right) - \exp \left(- \frac{2(t-1)(1-B/T_{\max})^2}{(1-\epsilon)^2} \right)$. Taking a

561 union bound gives the desired result.