

MipKV: A Sparsify-then-Recover Paradigm for Accelerating Large Vision-Language Model Pre-Filling

Supplementary Material

6. Detailed Evaluation Results

Detailed results for regular image-understanding tasks are summarized in Table 2, and high-resolution or video-related benchmarks are reported in Table 1.

7. More Visual Demonstration

Figures 13 to 17 show a complete attention distribution in all decoder layers. The attention sink pattern appears in most layers while intermittent scattering show in layer 14 and 20.

Figures 18 to 22 show the KV Deviation with and without

the token recovery in all decoder layers. High KV Deviations are largely suppressed.

8. Limitations

Due to the computational cost and resource constraints associated with evaluating very large multimodal models, our current experiments focus on Qwen2-VL-2B and LLaVA-1.5-7B. Extending MipKV to higher-capacity models—such as 72B-scale variants—and testing its generality across different LVLM architectures would provide a more comprehensive understanding of its applicability.

Table 1. Performance of MipKV compared with other acceleration methods on Qwen2-VL-2B across high-resolution and video benchmarks. Attention-based methods are excluded for their high memory consumption.

Methods	InfoVQA	DocVQA	HRBench ^{4K}	HRBench ^{8K}	MVBench	Avg.(%)
<i>upper bound, 100% visual tokens</i>						
Vanilla	66.0	90.1	59.2	58.1	66.0	100
<i>retain or recover 80% visual tokens</i>						
Resize	<u>64.0</u>	90.0	62.1	<u>56.2</u>	60.3	99.7
DART	<u>61.0</u>	88.0	58.2	<u>54.9</u>	59.3	96.3
V ² Drop	61.0	89.0	58.0	55.8	59.5	96.8
MipKV	65.0	90.0	<u>59.1</u>	57.9	<u>60.0</u>	<u>99.5</u>
<i>retain or recover 60% visual tokens</i>						
Resize	62.0	90.0	60.1	<u>54.8</u>	59.8	97.8
DART	54.0	80.0	58.4	<u>54.4</u>	58.5	92.0
V ² Drop	57.0	85.0	56.9	<u>53.6</u>	57.2	92.8
MipKV	<u>61.0</u>	<u>89.0</u>	<u>58.9</u>	57.4	<u>59.1</u>	<u>97.5</u>
<i>retain or recover 40% visual tokens</i>						
Resize	60.0	90.0	60.4	53.1	59.6	96.7
DART	43.0	66.0	57.1	<u>53.5</u>	55.9	84.0
V ² Drop	48.0	78.0	53.6	<u>52.9</u>	53.6	86.0
MipKV	<u>58.0</u>	<u>86.0</u>	<u>59.0</u>	56.8	<u>58.5</u>	<u>95.6</u>
<i>retain or recover 20% visual tokens</i>						
Resize	<u>52.0</u>	89.0	54.9	51.9	<u>56.1</u>	<u>90.5</u>
DART	<u>30.0</u>	43.0	<u>56.1</u>	<u>50.6</u>	<u>50.7</u>	<u>71.9</u>
V ² Drop	35.0	60.0	51.6	49.2	49.6	74.8
MipKV	55.0	<u>82.0</u>	58.0	56.0	57.5	92.8
<i>retain or recover 10% visual tokens</i>						
Resize	44.0	84.0	54.2	46.9	52.1	83.7
DART	24.0	26.0	<u>54.2</u>	<u>48.2</u>	45.8	<u>63.1</u>
V ² Drop	27.0	44.0	48.4	46.9	46.4	65.9
MipKV	52.0	<u>81.0</u>	57.4	55.6	56.9	91.2

Table 2. **Performance of MipKV compared with other acceleration methods on Qwen2-VL-2B across image-understanding benchmarks.** All experiments are conducted on a NVIDIA RTX 4090 GPU. The best results for each task are highlighted in bold, and underlined numbers indicate the second-best performance. “-” denotes out-of-memory errors, and “#” indicates that the input images are too small to be processed.

Methods	MME	MMB	SQA	GQA	POPE	ChartQA	OcrBench	VQA ^{text}	Avg.(%)
<i>upper bound, 100% visual tokens</i>									
Vanilla	1875.2	71.5	58.8	59.9	86.0	72.3	765	79.5	100
<i>retain or recover 80% visual tokens</i>									
Resize	1865.0	70.1	59.1	59.4	85.8	71.9	<u>736.0</u>	78.3	98.9
FastV	-	70.3	58.4	<u>59.6</u>	<u>85.7</u>	<u>68.3</u>	-	-	98.3
VisionZip	1834.8	70.2	58.5	<u>58.7</u>	85.6	64.1	639.0	70.0	94.2
DART	1916.0	69.4	58.2	59.4	85.5	69.1	697.0	<u>78.8</u>	97.8
V ² Drop	1878.3	<u>70.7</u>	57.9	59.5	85.5	68.4	669.0	78.2	97.1
MipKV	<u>1880.2</u>	71.1	<u>59.0</u>	59.9	86.0	72.3	761.0	79.1	99.9
<i>retain or recover 60% visual tokens</i>									
Resize	1867.6	69.1	60.0	58.5	84.9	68.3	<u>715.0</u>	<u>77.1</u>	97.4
FastV	-	68.4	58.9	58.6	84.0	55.7	-	-	93.7
VisionZip	1816.4	68.4	58.4	57.6	85.0	60.0	557.0	68.9	91.2
DART	1852.7	66.6	59.0	58.2	83.6	58.9	577.0	76.4	92.5
V ² Drop	1880.6	69.3	59.1	58.7	84.5	59.8	615.0	76.0	94.1
MipKV	1892.1	70.8	<u>59.3</u>	59.7	86.1	70.7	749.0	78.9	99.4
<i>retain or recover 40% visual tokens</i>									
Resize	<u>1832.3</u>	68.2	57.3	55.6	82.0	65.0	<u>680.0</u>	73.0	93.7
FastV	-	66.3	57.4	56.3	81.0	34.0	-	-	85.1
VisionZip	1800.4	65.8	56.9	56.1	84.3	54.5	460.0	65.9	86.9
DART	1779.6	63.1	<u>57.5</u>	55.8	79.7	47.5	455.0	71.1	85.2
V ² Drop	1824.6	68.2	<u>57.5</u>	56.5	82.4	43.2	525.0	<u>73.7</u>	87.7
MipKV	1883.5	70.8	60.0	59.5	86.6	68.2	739.0	78.3	98.9
<i>retain or recover 20% visual tokens</i>									
Resize	<u>1776.7</u>	63.9	55.4	52.3	76.2	40.7	#	66.1	84.8
FastV	-	58.2	54.1	50.8	71.1	18.5	-	-	73.3
VisionZip	1763.4	57.6	53.4	51.6	<u>79.5</u>	<u>41.2</u>	276.0	56.4	76.0
DART	1596.2	54.4	<u>57.4</u>	50.4	68.9	32.4	295.0	55.4	72.0
V ² Drop	1763.4	62.7	53.8	52.4	77.2	24.4	397.0	66.4	77.5
MipKV	1846.5	70.2	59.5	59.5	86.8	64.1	710.0	77.7	97.2
<i>retain or recover 10% visual tokens</i>									
Resize	<u>1710.0</u>	60.4	53.4	48.6	70.8	24.1	#	56.0	76.2
FastV	-	41.4	50.9	45.0	55.8	15.0	-	-	61.0
VisionZip	1503.9	44.9	53.7	46.2	<u>71.0</u>	<u>29.3</u>	140.0	41.1	63.1
DART	1449.8	47.0	<u>55.6</u>	46.0	55.6	21.3	190.0	38.1	60.2
V ² Drop	1605.2	57.2	53.2	48.6	69.1	16.0	<u>286.0</u>	59.0	68.9
MipKV	1855.1	69.7	59.2	59.3	86.8	60.0	688.0	77.0	95.8

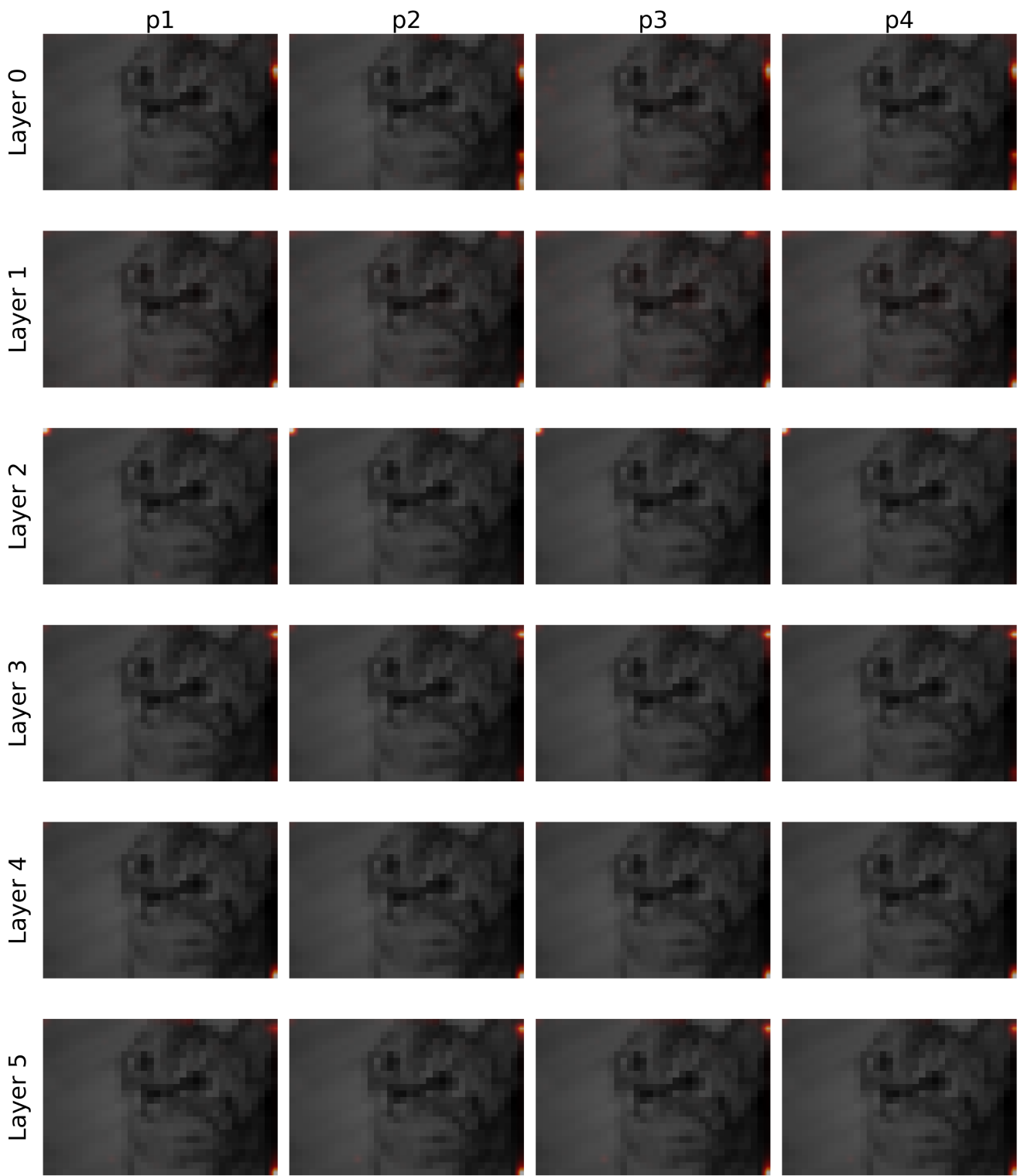


Figure 13. Attention distribution (Layer 0 – Layer 5)

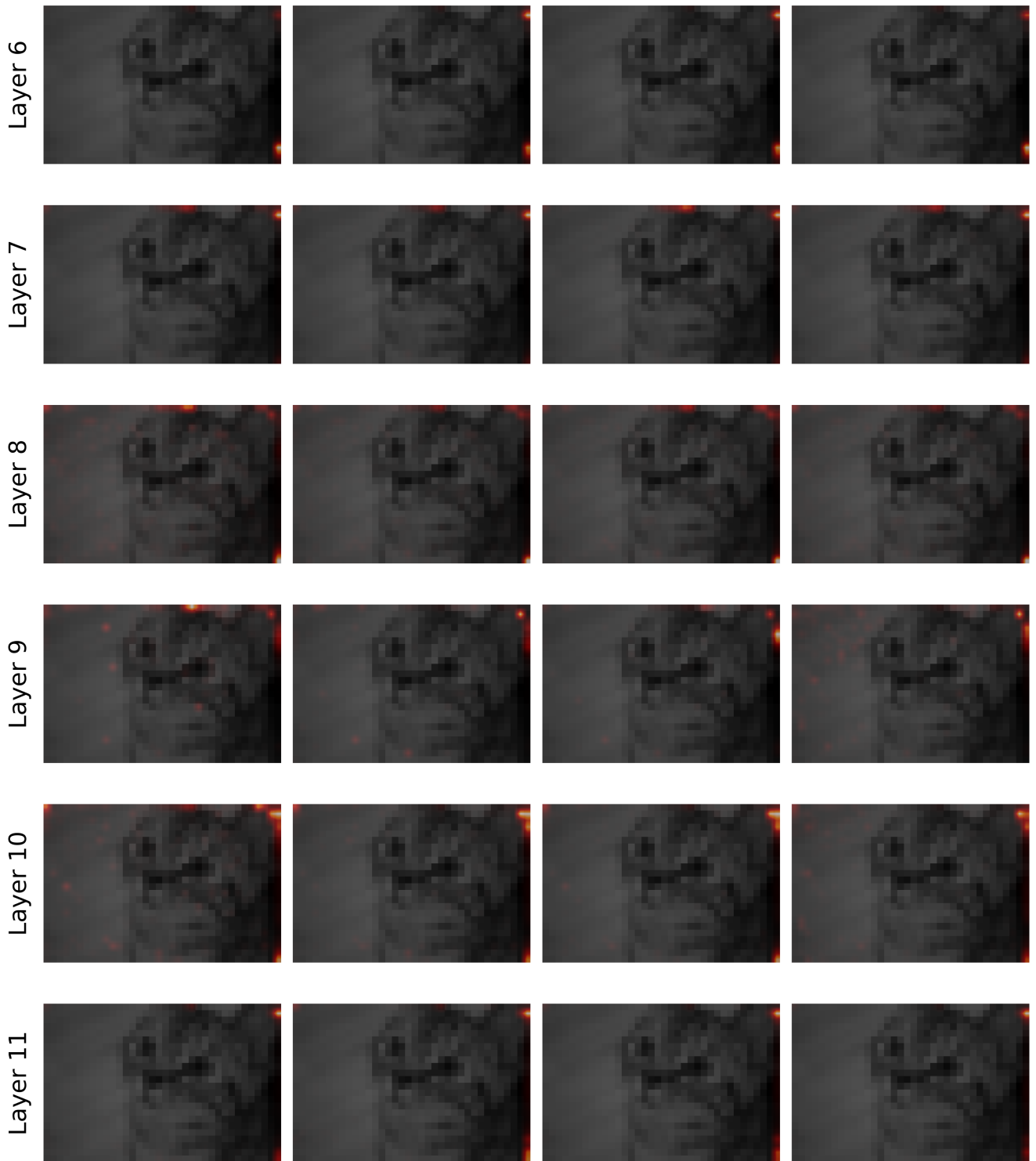


Figure 14. Attention distribution (Layer 6 – Layer 11)

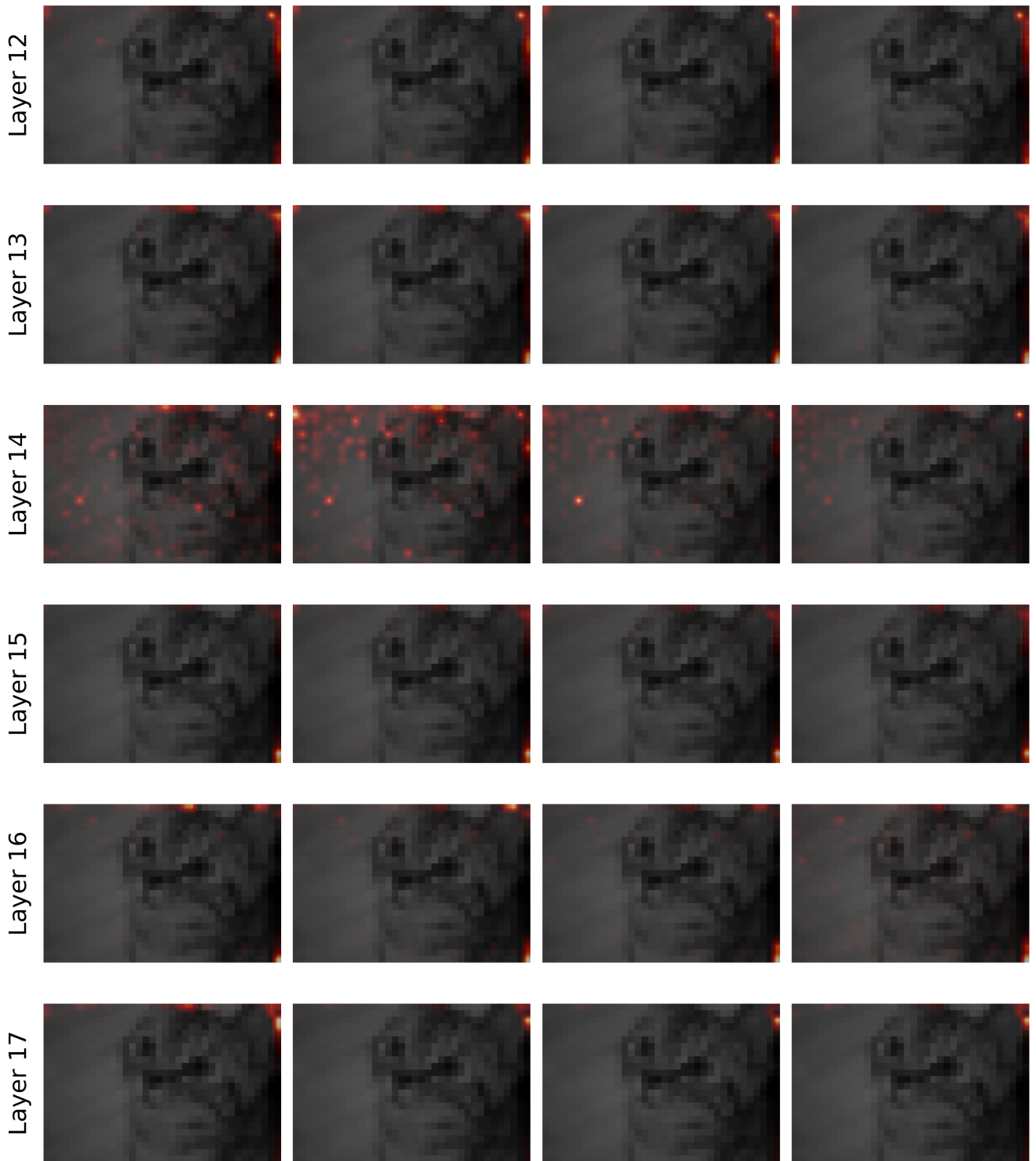


Figure 15. Attention distribution (Layer 12 – Layer 17)

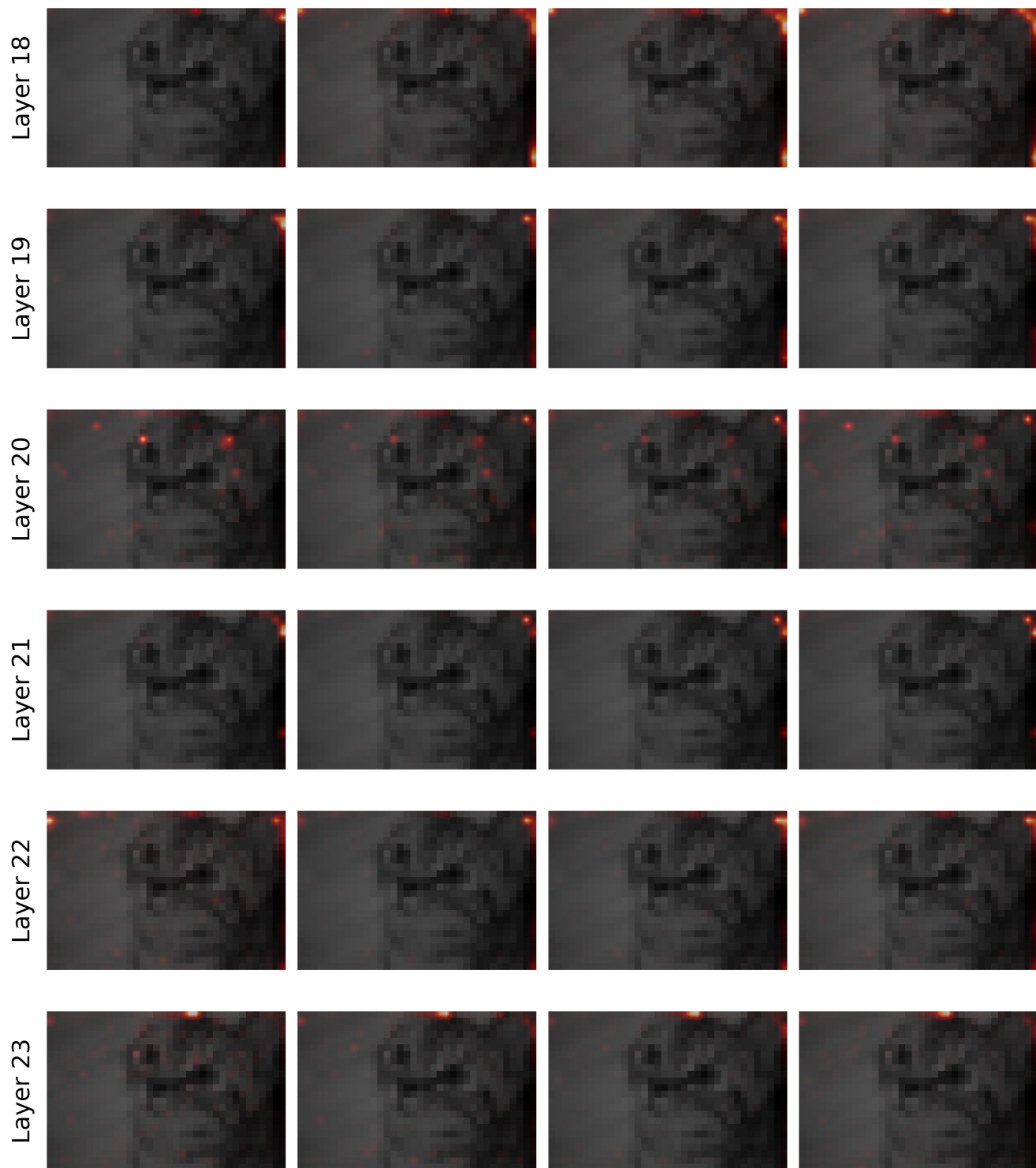


Figure 16. Attention distribution (Layer 18 – Layer 23)

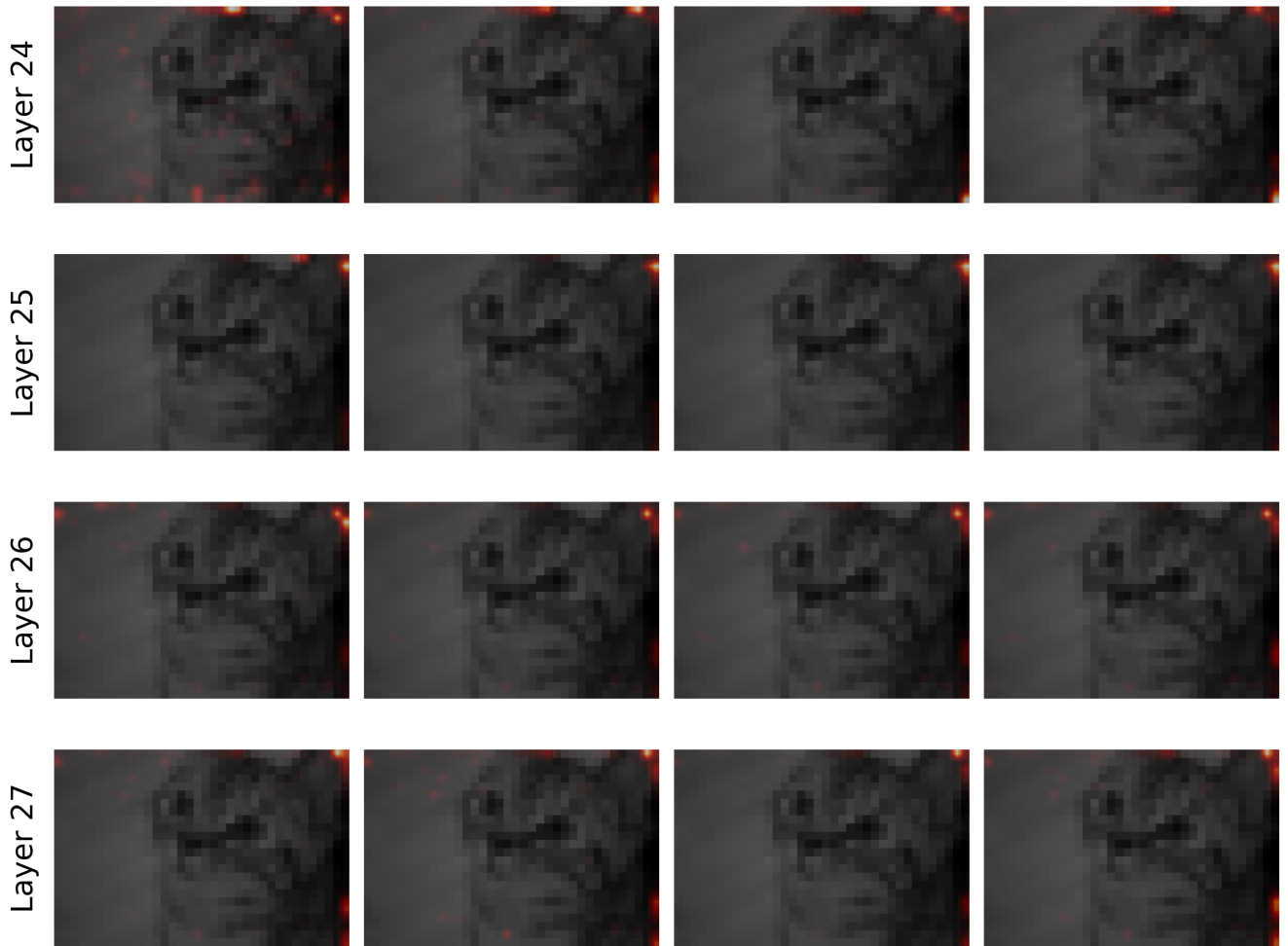


Figure 17. Attention distribution (Layer 24 – Layer 27)

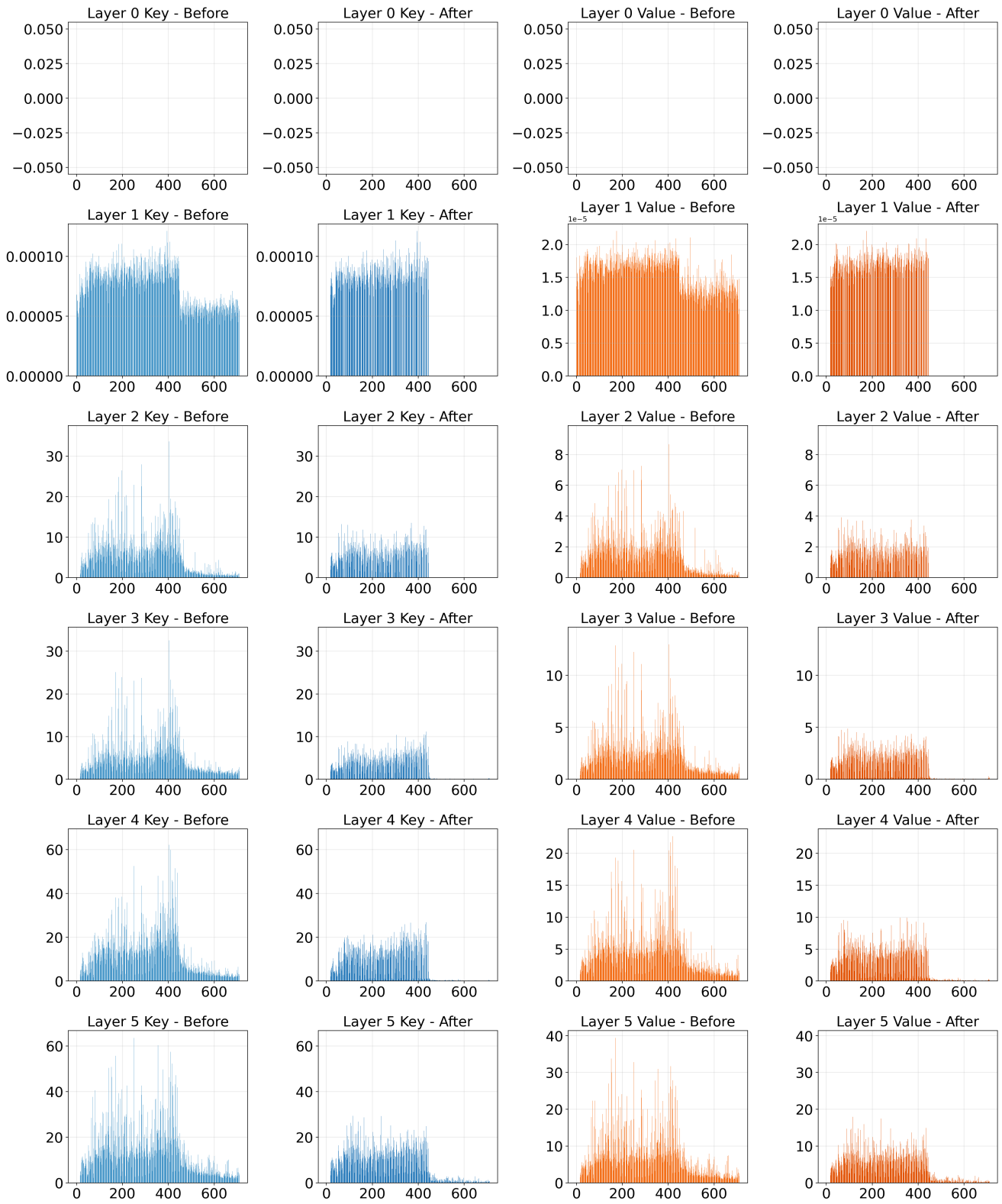


Figure 18. KV Deviation (Layer 0 – Layer 5)

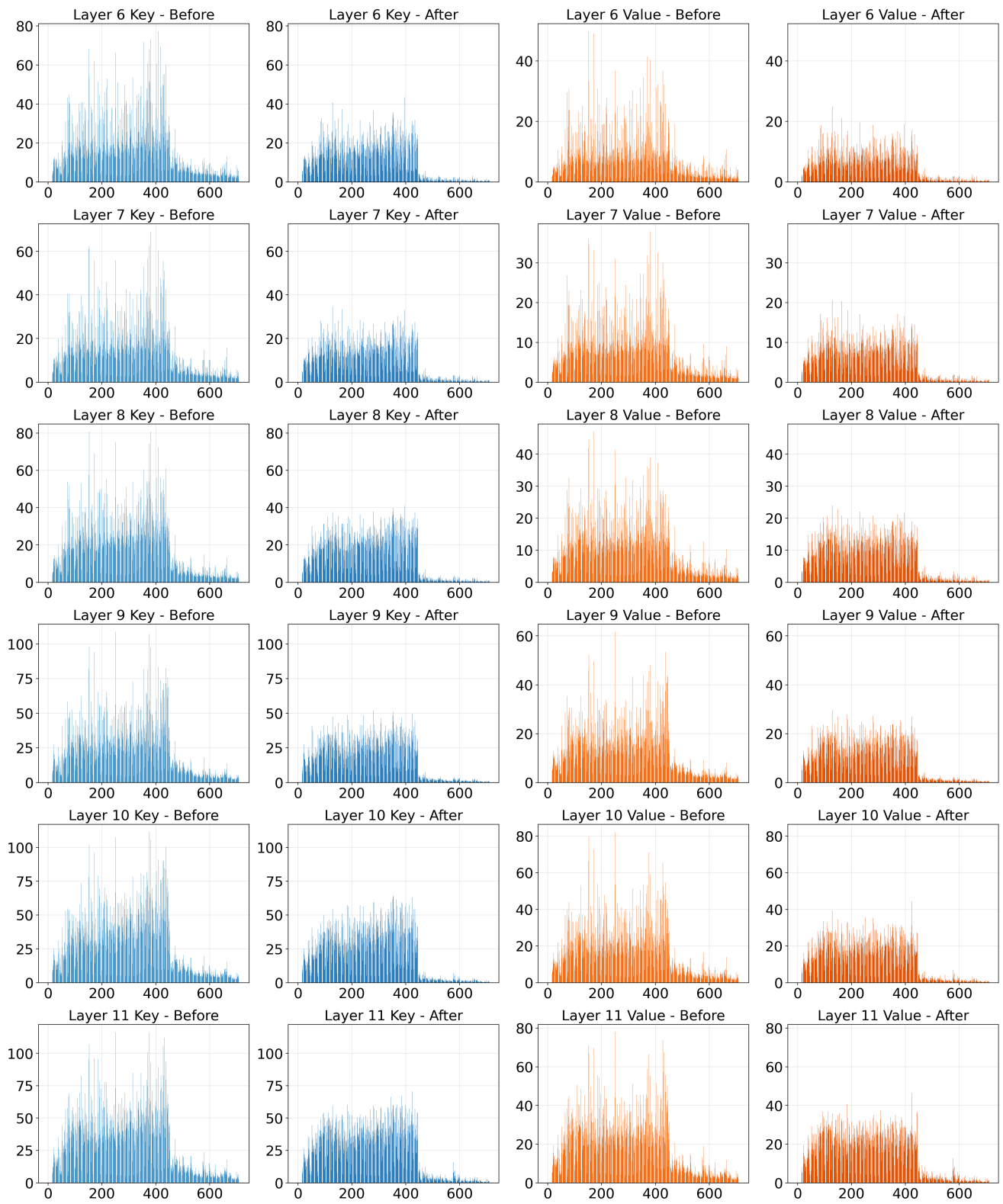


Figure 19. KV Deviation (Layer 6 – Layer 11)

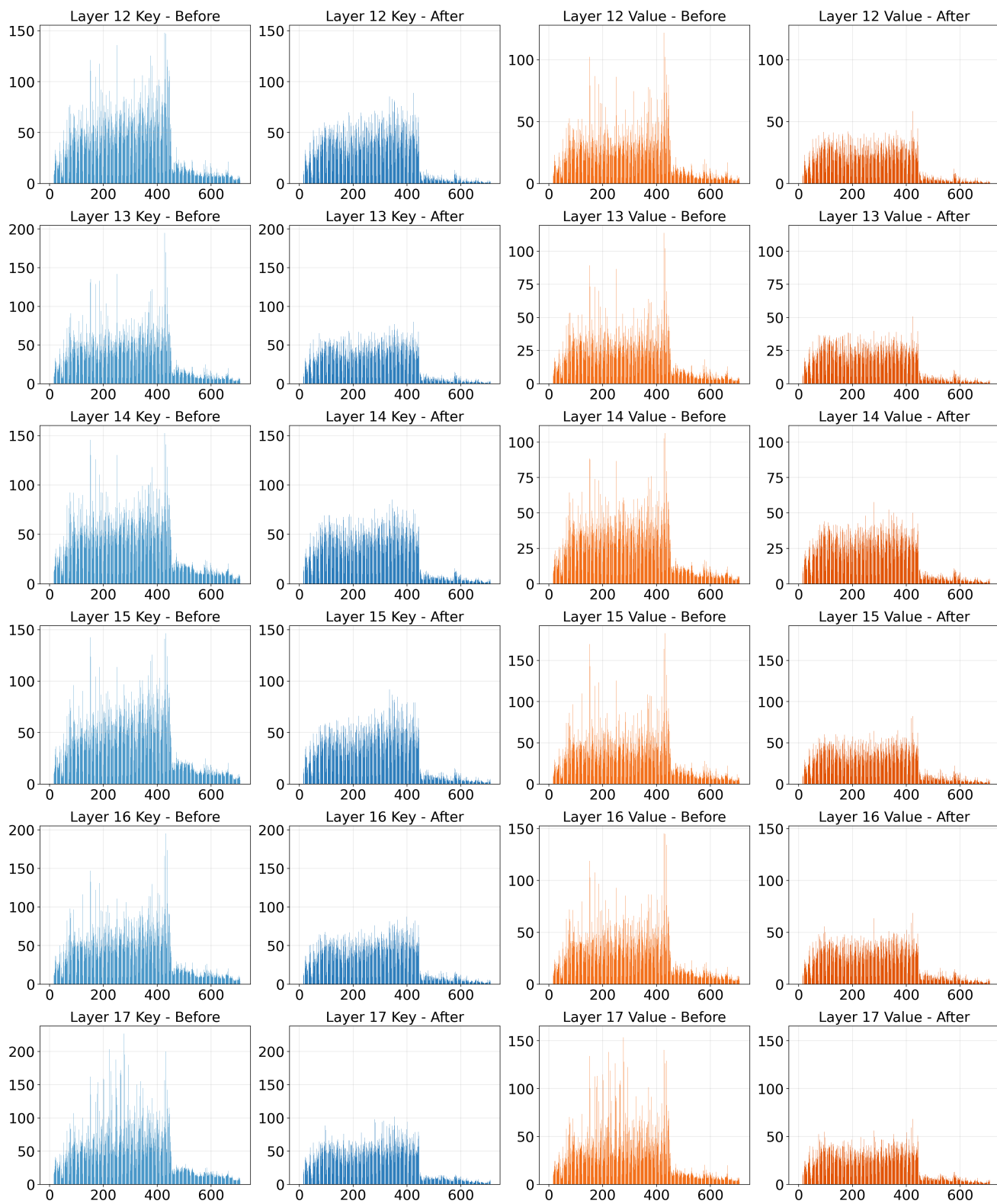


Figure 20. KV Deviation (Layer 12 – Layer 17)

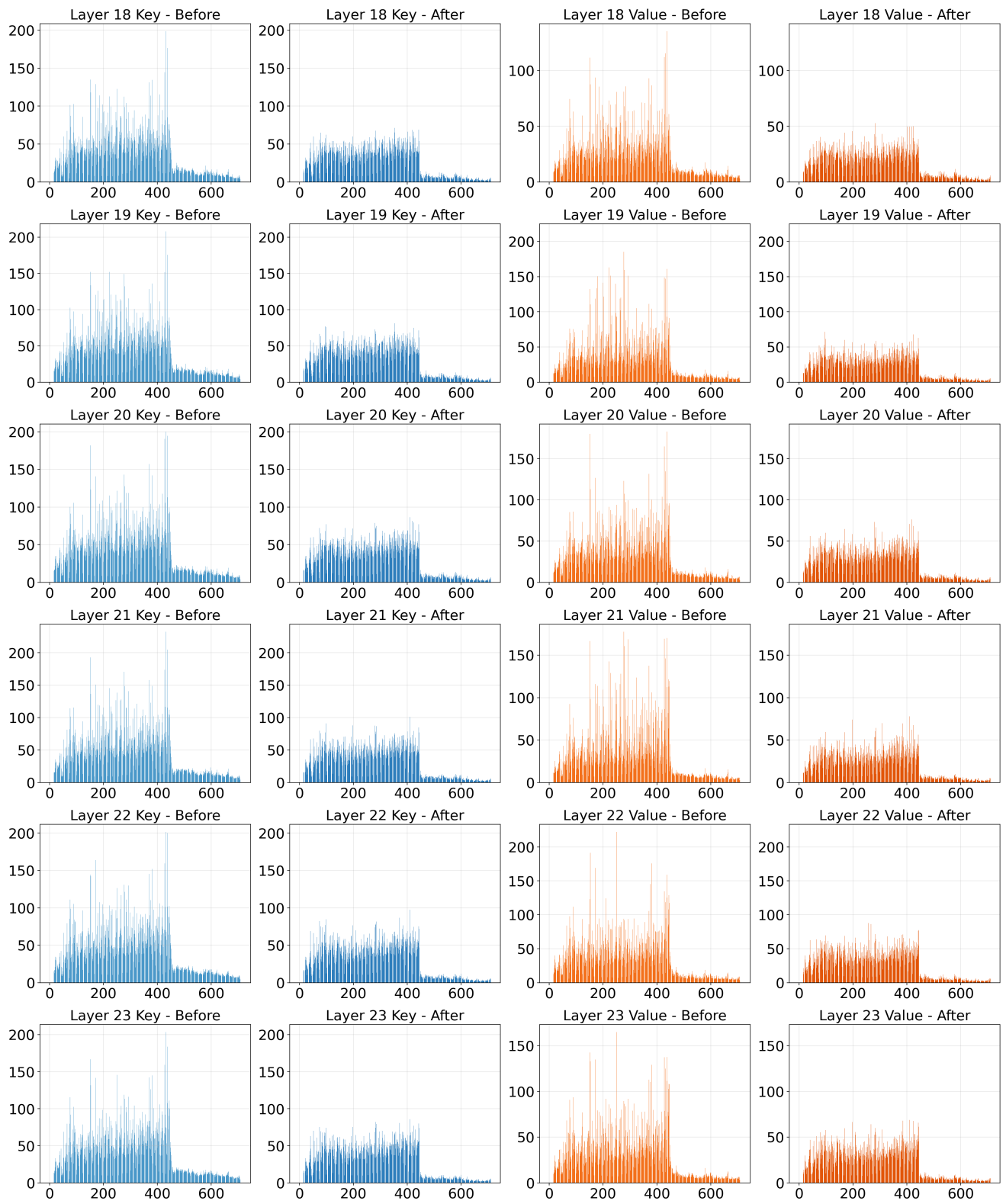


Figure 21. KV Deviation (Layer 18 – Layer 23)

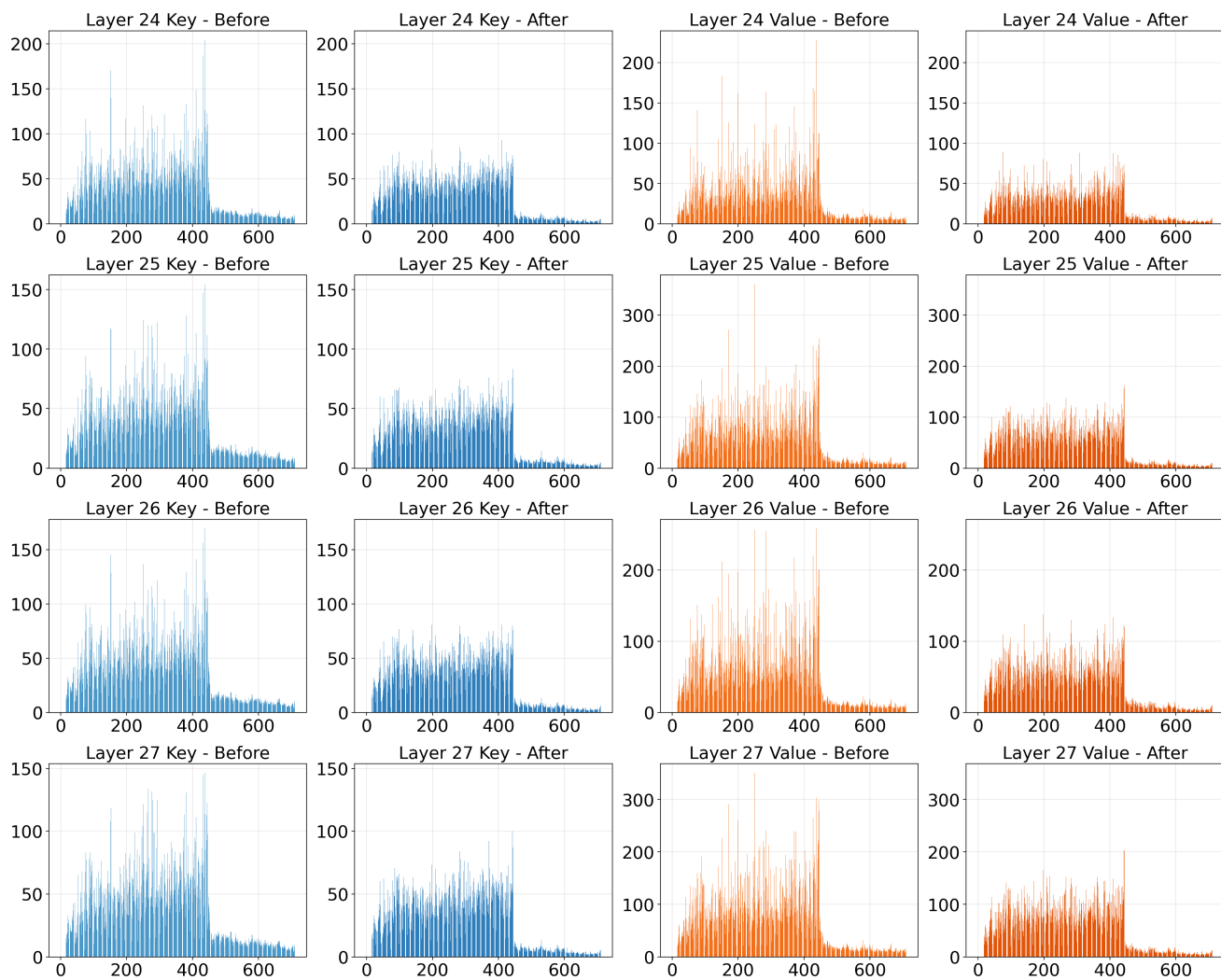


Figure 22. KV Deviation (Layer 24 – Layer 27)