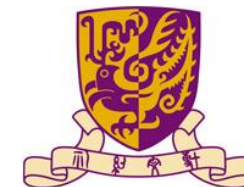# VideoRoPE

## What Makes for Good Video Rotary Position Embedding?

Xilin Wei*, Xiaoran Liu*, Yuhang Zang, Xiaoyi Dong,
Pan Zhang, Yuhang Cao, Jian Tong, Haodong Duan,
Qipeng Guo, Jiaqi Wang, Xipeng Qiu, Dahua Lin

Current advanced position embedding methods, such as M-RoPE, are still susceptible to periodic distractor interference (V-NIAH-D).

**haystack**   **needle**   **haystack**



Question: what is being transferred to the beaker in the laboratory?
A. Solid substance   B. Gas   C. Nothing   D. Liquid tester
M-RoPE: A. Solid substance 😅
VideoRoPE: D. Liquid tester 😃



M-RoPE - t   M-RoPE - x   M-RoPE - y
VideoRoPE - t   VideoRoPE - x   VideoRoPE - y

- Both M-RoPE and VideoRoPE **successfully locate the needle information** required to answer the question.

- Due to suboptimal frequency allocation, M-RoPE focuses on vertical cues at the expense of temporal semantics, leading to **poor long-range modeling and wrong answers.** VideoRoPE, by leveraging **temporal localization, answers correctly.**

3

- most previous methods only cover **part of** the table

- **VideoRoPE achieves a full-stack design,** addressing **all four core dimensions**: structural modeling, frequency allocation, spatial symmetry, and temporal scaling—surpassing prior RoPE variants.
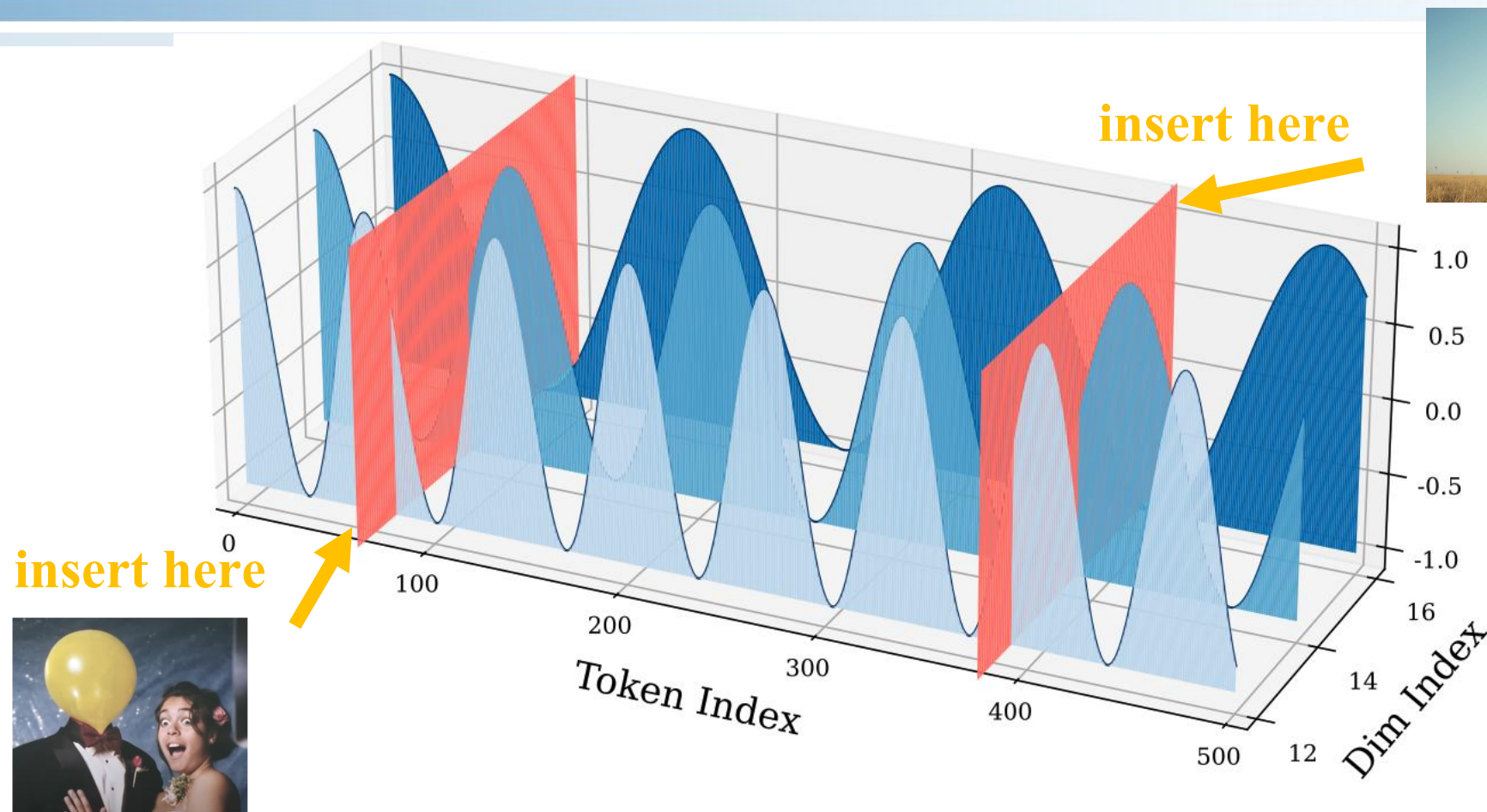
| | 2D/3D Structure | Frequency Allocation | Spatial Symmetry | Temporal Index Scaling | al etry | Temporal Index Scaling |
|---|---|---|---|---|---|---|
| Vanilla RoPE (Su et al., 2024) | ✗ | | ✗ | ✗ | | ✗ |
| TAD-RoPE (Gao et al., 2024) | ✗ | | ✗ | ✗ | | ✓ |
| RoPE-Tie (Su, 2024a) | ✓ | | ✗ | ✓ | | ✗ |
| M-RoPE (Wang et al., 2024a) | ✓ | | ✗ | ✗ | | ✗ |
| VideoRoPE (Ours) | ✓ | | ✓ | ✓ | | ✓ |

(a) Temporal Frequency Allocation in M-RoPE

As shown in the red planes, positions that are far apart in time can end up with **similar** positional encodings due to these **oscillations**.

compare



(a) Temporal Frequency Allocation in M-RoPE

(b) Temporal Frequency Allocation in VideoRoPE (ours)

VideoRoPE adopts **low-frequency modeling for the temporal dimension**, achieving better long-range monotonicity and **avoiding oscillations**, which effectively reduces distractor interference in V-NIAH-D.

(a) Temporal Frequency Allocation in M-RoPE

$$\left(\begin{array}{c} q^{(96)} \\ q^{(97)} \\ q^{(98)} \\ q^{(99)} \\ \vdots \\ q^{(126)} \\ q^{(127)} \end{array}\right)^{\top} \left(\begin{array}{ccccccc} \cos\theta_{48}\Delta t & -\sin\theta_{48}\Delta t & 0 & 0 & \cdots & 0 & 0 \\ \sin\theta_{48}\Delta t & \cos\theta_{48}\Delta t & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos\theta_{49}\Delta t & 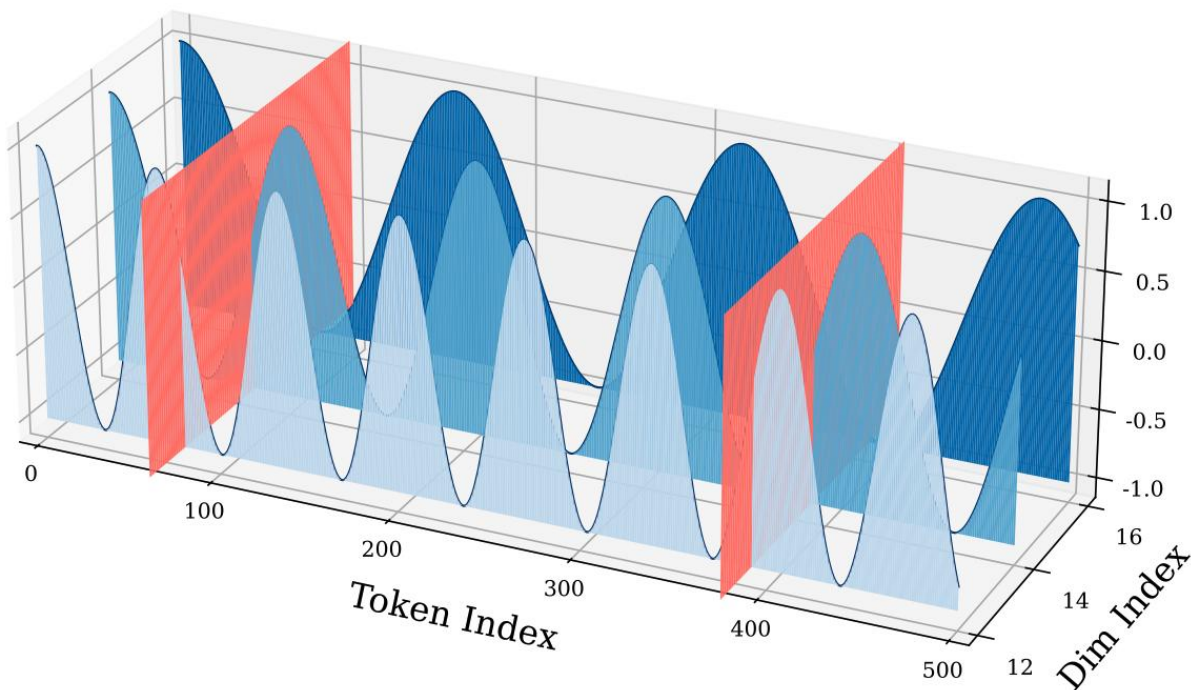-\sin\theta_{49}\Delta t & \cdots & 0 & 0 \\ 0 & 0 & \sin\theta_{49}\Delta t & \cos\theta_{49}\Delta t & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos\theta_{63}\Delta t & -\sin\theta_{63}\Delta t \\ 0 & 0 & 0 & 0 & \cdots & \sin\theta_{63}\Delta t & \cos\theta_{63}\Delta t \end{array}\right) \left(\begin{array}{c} k^{(96)} \\ k^{(97)} \\ k^{(98)} \\ k^{(99)} \\ \vdots \\ k^{(126)} \\ k^{(127)} \end{array}\right)$$
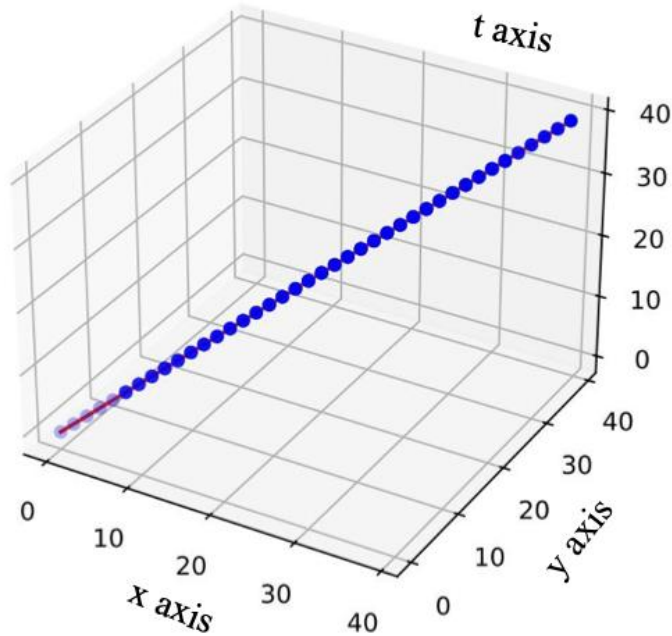
modeling temporal dependency with lower frequency ☺

$$+ \left(\begin{array}{c} q^{(0)} \\ q^{(1)} \\ q^{(4)} \\ q^{(5)} \\ \vdots \\ q^{(92)} \\ q^{(93)} \end{array}\right)^{\top} \left(\begin{array}{ccccccc} \cos\theta_{0}\Delta x & -\sin\theta_{0}\Delta x & 0 & 0 & \cdots & 0 & 0 \\ \sin\theta_{0}\Delta x & \cos\theta_{0}\Delta x & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos\theta_{2}\Delta x & -\sin\theta_{2}\Delta x & \cdots & 0 & 0 \\ 0 & 0 & \sin\theta_{2}\Delta x & \cos\theta_{2}\Delta x & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos\theta_{46}\Delta x & -\sin\theta_{46}\Delta x \\ 0 & 0 & 0 & 0 & \cdots & \sin\theta_{46}\Delta x & \cos\theta_{46}\Delta x \end{array}\right) \left(\begin{array}{c} k^{(0)} \\ k^{(1)} \\ k^{(4)} \\ k^{(5)} \\ \vdots \\ k^{(92)} \\ k^{(93)} \end{array}\right)$$
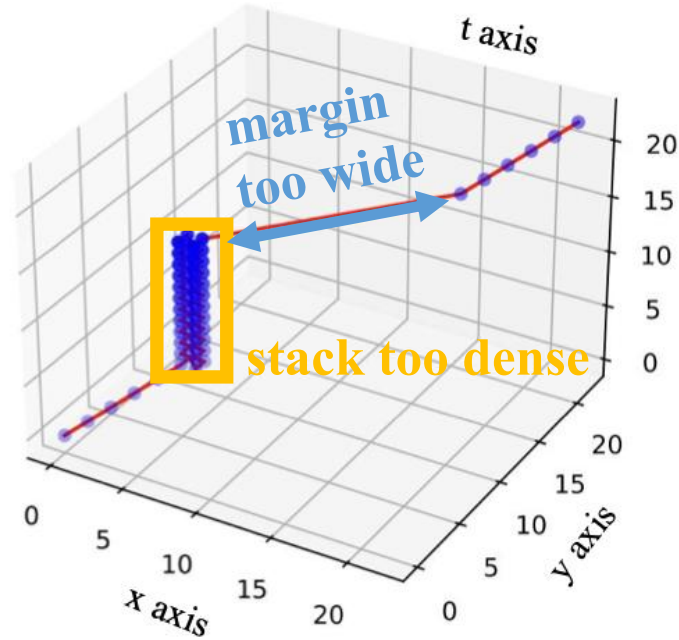
modeling horizontal dependency with interleaved high frequency

$$+ \left(\begin{array}{c} q^{(2)} \\ q^{(3)} \\ q^{(6)} \\ q^{(7)} \\ \vdots \\ q^{(94)} \\ q^{(95)} \end{array}\right)^{\top} \left(\begin{array}{ccccccc} \cos\theta_{1}\Delta y & -\sin\theta_{1}\Delta y & 0 & 0 & \cdots & 0 & 0 \\ \sin\theta_{1}\Delta y & \cos\theta_{1}\Delta y & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos\theta_{3}\Delta y & -\sin\theta_{3}\Delta y & \cdots & 0 & 0 \\ 0 & 0 & \sin\theta_{3}\Delta y & \cos\theta_{3}\Delta y & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos\theta_{47}\Delta y & -\sin\theta_{47}\Delta y \\ 0 & 0 & 0 & 0 & \cdots & \sin\theta_{47}\Delta y & \cos\theta_{47}\Delta y \end{array}\right) \left(\begin{array}{c} k^{(2)} \\ k^{(3)} \\ k^{(6)} \\ k^{(7)} \\ \vdots \\ k^{(94)} \\ k^{(95)} \end{array}\right)$$

modeling vertical dependency with interleaved high frequency
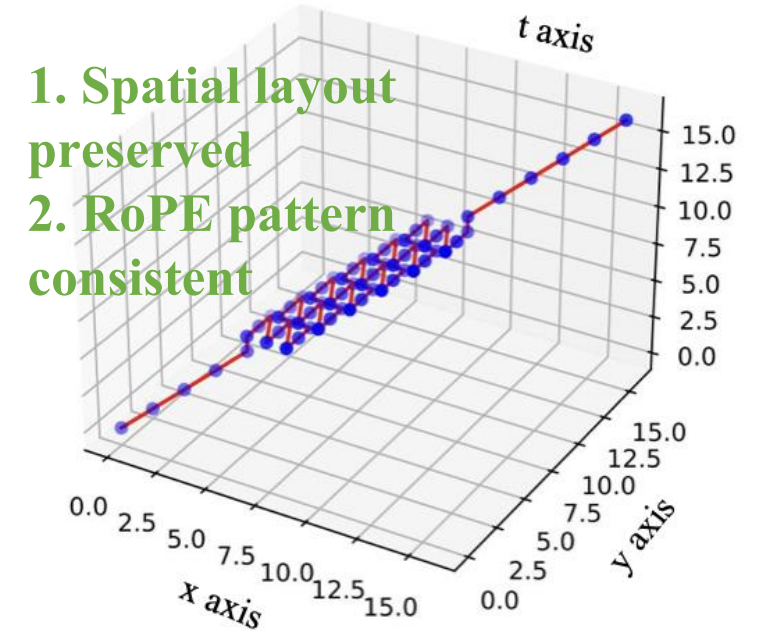
(a) 3D visualization for Vanilla RoPE.

(b) 3D visualization for M-RoPE.

(c) 3D visualization for VideoRoPE.

3D visualization of position embeddings: (a) Vanilla RoPE **lacks spatial modeling**. (b) M-RoPE introduces **inconsistent index growth** across frames. (c) VideoRoPE **balances spatial modeling with consistent indexing**, preserving RoPE's desirable structure.

8

$$
(t, x, y) = \begin{cases}
(\tau, \tau, \tau) & \text{if } 0 \leq \tau < T_s \\[2ex]
\begin{pmatrix} T_s + \delta(\tau - T_s), \\ T_s + \delta(\tau - T_s) + w - \frac{W}{2}, \\ T_s + \delta(\tau - T_s) + h - \frac{H}{2} \end{pmatrix} & \text{if } T_s \leq \tau < T_s + T_v \\[4ex]
\begin{pmatrix} \tau + (\delta - 1)T_v, \\ \tau + (\delta - 1)T_v, \\ \tau + (\delta - 1)T_v \end{pmatrix} & \text{if } T_s + T_v \leq \tau < T_s + T_v + T_e
\end{cases}
$$

Adjustable Temporal Spacing(ATS). To scale the temporal index, we introduce a scaling factor δ to better align temporal information between visual and textual tokens.

$$(t, x, y) = \begin{cases} (\tau, \tau, \tau) & \text{if } 0 \leq \tau < T_s \\ \begin{pmatrix} T_s + \delta(\tau - T_s), \\ T_s + \delta(\tau - T_s) + w - \frac{W}{2}, \\ T_s + \delta(\tau - T_s) + h - \frac{H}{2} \end{pmatrix} & \text{if } T_s \leq \tau < T_s + T_v \\ \begin{pmatrix} \tau + (\delta - 1)T_v, \\ \tau + (\delta - 1)T_v, \\ \tau + (\delta - 1)T_v \end{pmatrix} & \text{if } T_s + T_v \leq \tau < T_s + T_v + T_e \end{cases}$$

Adjustable Temporal Spacing(ATS). To scale the temporal index, we introduce a scaling factor $\delta$ to better align temporal information between visual and textual tokens.
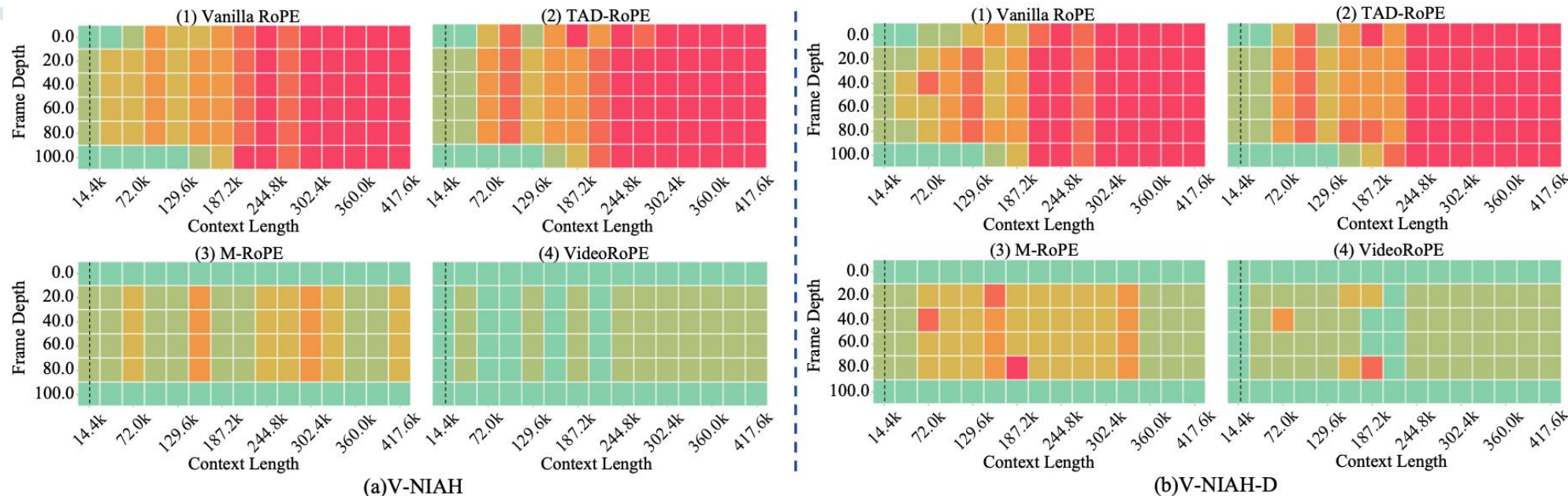
$$(t, x, y) = \begin{cases} (\tau, \tau, \tau) & \text{if } 0 \leq \tau < T_s \\ \begin{pmatrix} T_s + \delta(\tau - T_s), \\ T_s + \delta(\tau - T_s) + w - \frac{W}{2}, \\ T_s + \delta(\tau - T_s) + h - \frac{H}{2} \end{pmatrix} & \text{if } T_s \leq \tau < T_s + T_v \\ \begin{pmatrix} \tau + (\delta - 1)T_v, \\ \tau + (\delta - 1)T_v, \\ \tau + (\delta - 1)T_v \end{pmatrix} & \text{if } T_s + T_v \leq \tau < T_s + T_v + T_e \end{cases}$$

Adjustable Temporal Spacing(ATS). To scale the temporal index, we introduce a scaling factor $\delta$ to better align temporal information between visual and textual tokens.

# Experiments on Long Video Understanding

| Method | LongVideoBench | | | | MLVU | | | | Video-MME | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 8k | 16k | 32k | 64k | 8k | 16k | 32k | 64k | 8k | 16k | 32k | 64k |
| Vanilla RoPE (Su et al., 2024) | **54.97** | 54.87 | <u>54.56</u> | 54.04 | 63.31 | <u>65.79</u> | <u>65.93</u> | 62.02 | <u>60.67</u> | 60.00 | 61.33 | 58.33 |
| TAD-RoPE (Gao et al., 2024) | 54.14 | <u>55.08</u> | 53.94 | 53.42 | <u>63.67</u> | 65.28 | 65.28 | 60.73 | 60.33 | **61.33** | **62.00** | 58.67 |
| M-RoPE (Wang et al., 2024a) | 53.42 | 52.80 | 53.11 | 54.35 | 60.41 | 60.68 | 61.56 | 61.10 | <u>60.67</u> | 59.67 | 61.00 | 59.67 |
| VideoRoPE (Ours) | <u>54.46</u> | **55.29** | **57.15** | **57.26** | **65.19** | **66.29** | **66.02** | **65.56** | **61.33** | <u>61.00</u> | <u>61.67</u> | **61.33** |

☐ Benchmarks: LongVideoBench, MLVU, VideoMME

☐ Consistent gains over M-RoPE: +2.91 / +4.46 / +1.66 @64k context

☐ Robust to long-range dependencies

☐ Strong adaptability across tasks

(a)V-NIAH

(b)V-NIAH-D

| Method | V-NIAH Acc. | V-NIAH-D Acc. |
|---|---|---|
| Vanilla RoPE (Su et al., 2024) | 31.78 | 30.22 |
| TAD-RoPE (Gao et al., 2024) | 29.33 | 29.56 |
| M-RoPE (Wang et al., 2024a) | 78.67 | 74.67 |
| VideoRoPE | **91.11** | **87.11** |

- ☐ V-NIAH-D is more challenging than V-NIAH
- ☐ Vanilla RoPE / TAD-RoPE: limited extrapolation
- ☐ VideoRoPE > M-RoPE in long context extrapolation
- ☐ +12.44% over M-RoPE on Video Retrieval

13

| Method | OR | T | SD | F | NF | Avg. |
|---|---|---|---|---|---|---|
| Vanilla RoPE (Su et al., 2024) | 51.5 | 30.0 | 48.0 | 8.0 | 43.0 | 36.1 |
| TAD-RoPE (Gao et al., 2024) | 51.0 | 37.0 | 48.0 | 11.5 | 47.5 | 39.0 |
| M-RoPE (Wang et al., 2024a) | 39.0 | 29.0 | 43.5 | 12.5 | 47.5 | 34.3 |
| VideoRoPE | **57.0** | **58.5** | **50.5** | **15.0** | **50.0** | **46.2** |

- ❑ +29.5% on Temporal Hallucination → better temporal reasoning
- ❑ +18.0% on Spatial/Object-Relation Hallucination → better spatial understanding
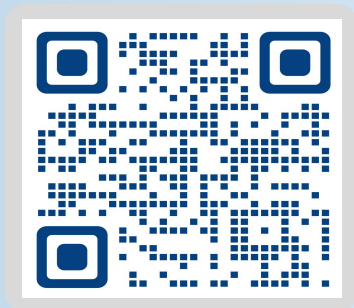- ❑ Robust to complex video hallucinations

| Method | LongVideoBench | | | | MLVU | | | |
|---|---|---|---|---|---|---|---|---|
| | 8k | 16k | 32k | 64k | 8k | 16k | 32k | 64k |
| Baseline | 53.42 | 52.80 | 53.11 | 54.35 | 60.41 | 60.68 | 61.56 | 61.10 |
| + DL | 52.17 | 52.07 | 53.31 | 53.63 | 62.06 | 63.03 | 62.52 | 62.75 |
| + DL & LTA | **54.46** | **55.49** | 54.66 | 55.60 | 63.35 | 64.09 | 64.00 | 63.26 |
| + DL & LTA & ATS | **54.46** | 55.29 | **57.15** | **57.26** | **65.19** | **66.29** | **66.02** | **65.56** |

- ☐ Ablation on LongVideoBench & MLVU (64k context)
- ☐ Baseline (M-RoPE): 54.35 / 61.10
- ☐ +DL → +LTA → +ATS → performance improves progressively
- ☐ Final: 57.26 / 65.56
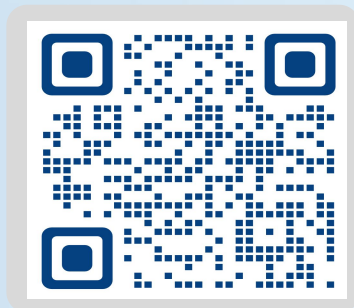- ☐ Effective use of spatio-temporal position encoding

\* Further ablations on layout strategies, frequency allocation, ATS scaling, and DL are provided in the main paper.

15

# Conclusion

☐ **Four key criteria** for effective positional encoding:
  ☐ 2D/3D structure, frequency allocation, spatial symmetry, temporal index scaling

☐ Prior RoPE variants struggle with **temporal distractors due to improper allocation**
  ☐ VideoRoPE **addresses** this with:
  ☐ 3D spatiotemporal structure
  ☐ Low-frequency temporal allocation (reduces oscillations)
  ☐ Diagonal spatial layout (ensures symmetry)
  ☐ Adjustable temporal spacing (ATS)

☐ **Superior performance** in:
  ☐ Long video retrieval
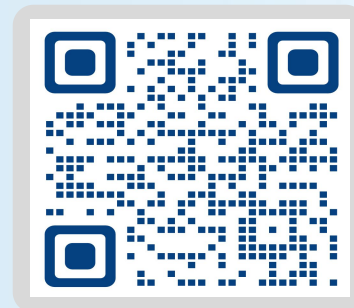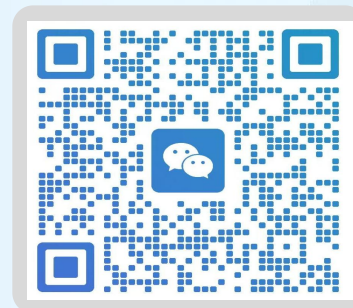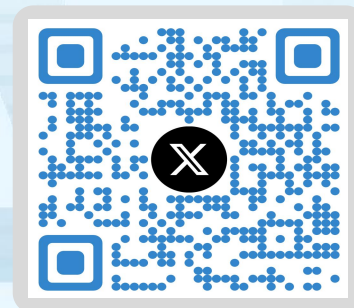  ☐ Video understanding
  ☐ Video hallucination tasks

| Thanks

Github     Paper     Homepage     wechat     X

Contact us: wiselnn570@gmail.com