



Xindi Wu^{1,2} Sven Elflein^{1,3,4} James Lucas¹ Olga Russakovsky²
 Laura Leal-Taixé¹ Despoina Paschalidou¹ Jonathan Lorraine¹ Aljosa Osep¹
¹NVIDIA · ²Princeton University · ³University of Toronto · ⁴Vector Institute

Out of Sight, Out of Mind

Autoregressive video world models promise interactive worlds - but **visual persistence** collapses once generation exceeds the training horizon.

Can an autoregressive video world model reliably remember where it has been, at any generation length?

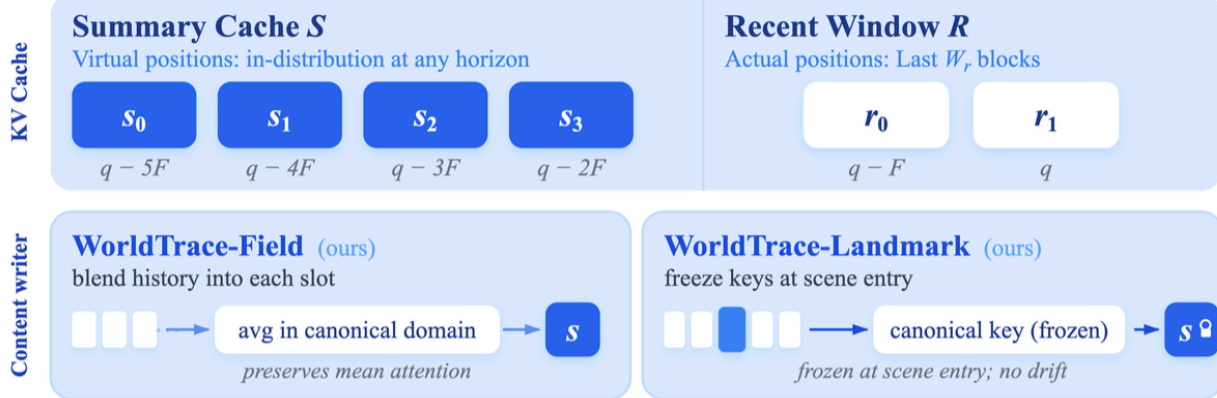
Position and Content are Two Coupled Bottlenecks

⊖ **Addressability:** Temporal RoPE offsets exceed the training range, making cached memories unreadable even when physically present.

⊖ **Content Fidelity:** Naive key averaging in RoPE-rotated space causes phase cancellation, destroying the signal that compressed summaries should carry.

Method: WorldTrace

WorldTrace



WorldTrace-Field

Averages *all* history into each summary slot in the **canonical (unrotated) key domain**, then re-rotates to the slot's virtual position.

$$K_{\text{field}}^{(k)}(t_v) = R(\theta_k t_v) \frac{1}{M} \sum_{m=1}^M R(-\theta_k t_m) K_{t_m}^{(k)}$$

Avoids phase cancellation by aligning keys to a shared canonical phase. **Goal:** coherence: drift-free generation at long horizons.

WorldTrace-Landmark

Detects scene-entry events via cosine-distance spikes in canonical-K space. Fills summary slots with **verbatim frames** at detected boundaries, frozen in canonical form.

$$K_{\text{land}}^{(k)}(t_v) = R(\theta_k t_v) R(-\theta_k t_{\ell^*}) K_{t_{\ell^*}}^{(k)}$$

Key stored once at landmark time; single rotation per shift; no float error accumulation. **Goal:** episodic recall after a long detour.

Slot-Rank Position Assignment

Virtual position for summary slot s :

$$v_s = q - (L_{\text{train}} - 1 - s) \cdot F$$

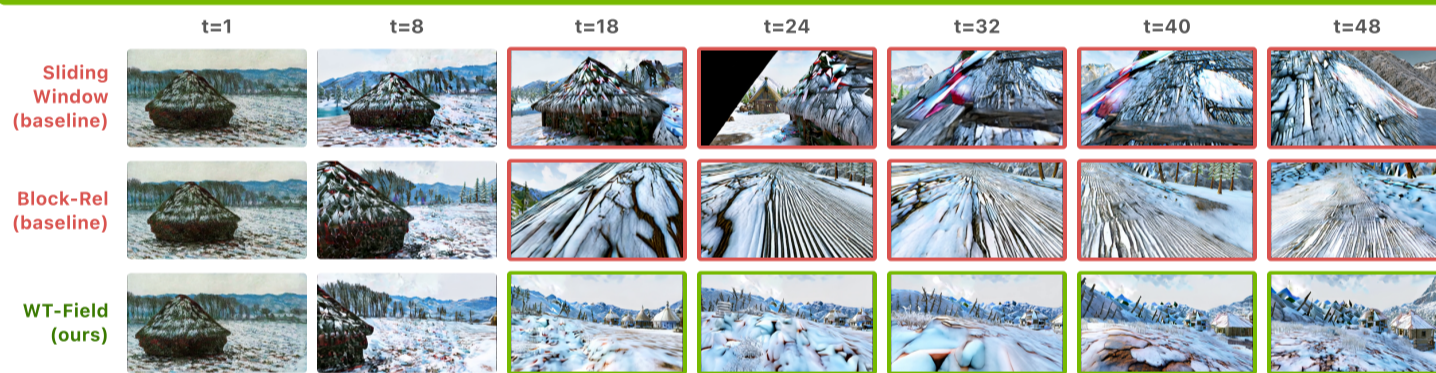
q = current query · L_{train} = training context length · F = frames per AR block

In-distribution at any generation length

Horizon-stable: rank-only offsets

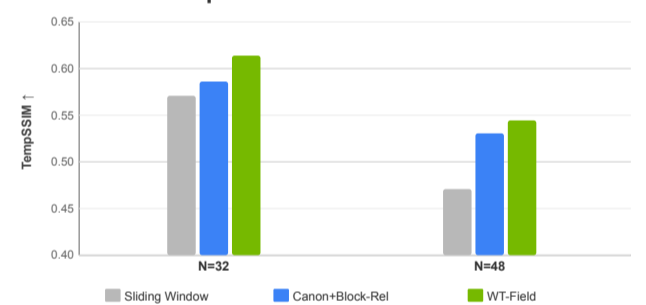
Experiments

WorldTrace-Field



Both baselines (Sliding Window, Block-Rel) diverge from $t=18$ ✗; WT-Field remains coherent through $t=48$ ✓.

Coherence: TempSSIM ↑

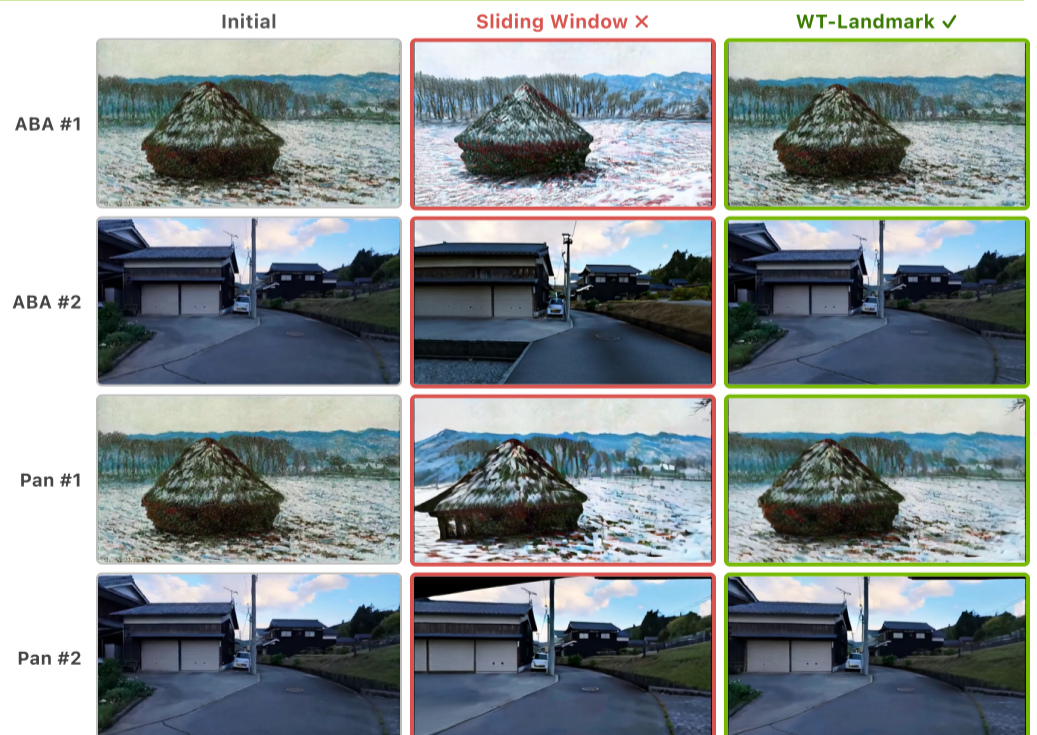
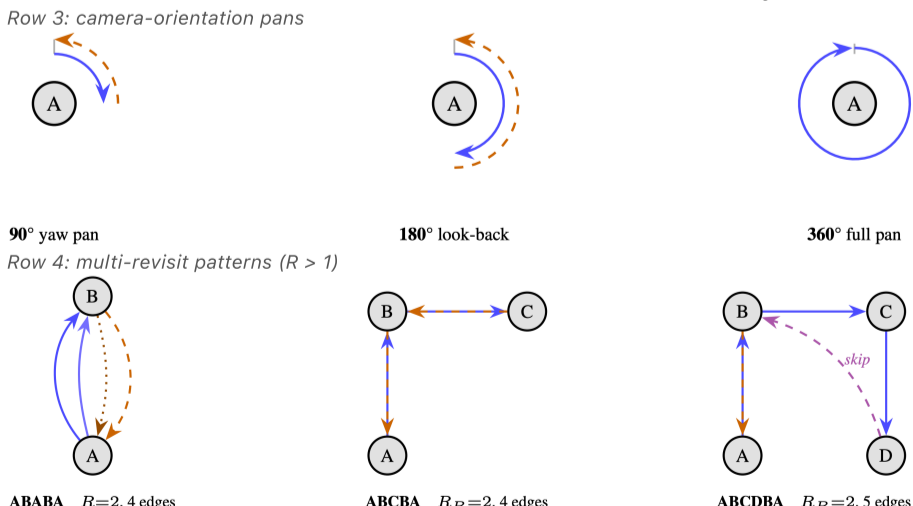
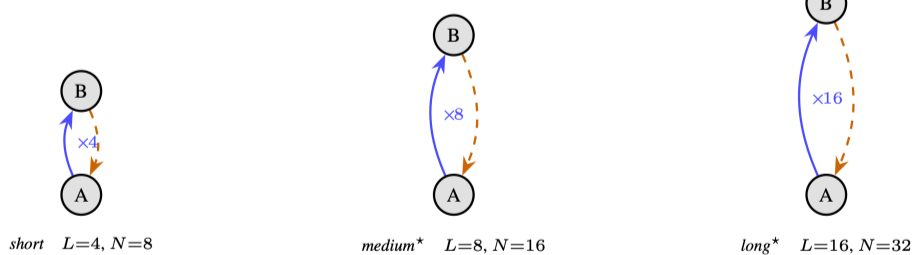
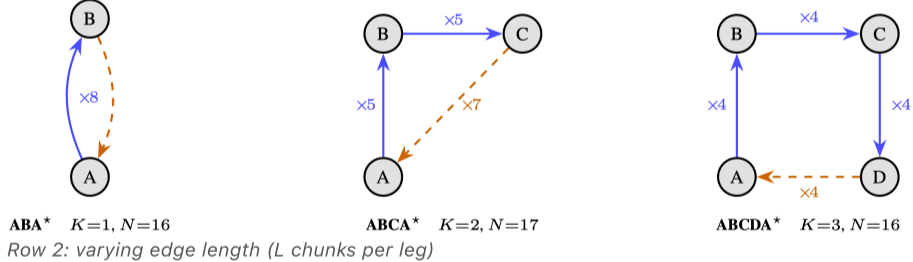


+15.5% relative TempSSIM at $N=48$; WT-Field leads on both horizons.

WorldTrace-Landmark

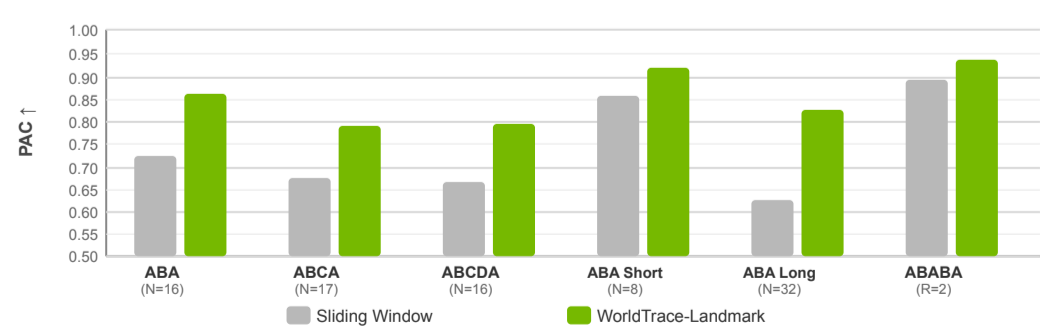
LoopMem Benchmark Topologies

Row 1: varying topology (K intermediate waypoints)



WT-Landmark faithfully restores scene A across diverse environments.

LoopMem PAC ↑ — Sliding Window vs. WT-Landmark



Gap widens with context distance (+6.3% to +19.8%). Consistent gains across revisit depth.

Conclusion

Long-horizon failure is fundamentally a problem of addressability.

OOD Positions Are the Binding Constraint

Content compression alone cannot restore recall once temporal RoPE offsets go out-of-distribution.

WT-Field +15.5% TempSSIM

Canonical key averaging improves long-range coherence at 24x generation horizon.

Training-Free

Only the KV cache changes: no fine-tuning, no model modification. Applies to any temporal-RoPE AR video model.

WT-Landmark +19.3–20.0% PAC

Frozen canonical landmark keys restore episodic recall across all LoopMem topologies.