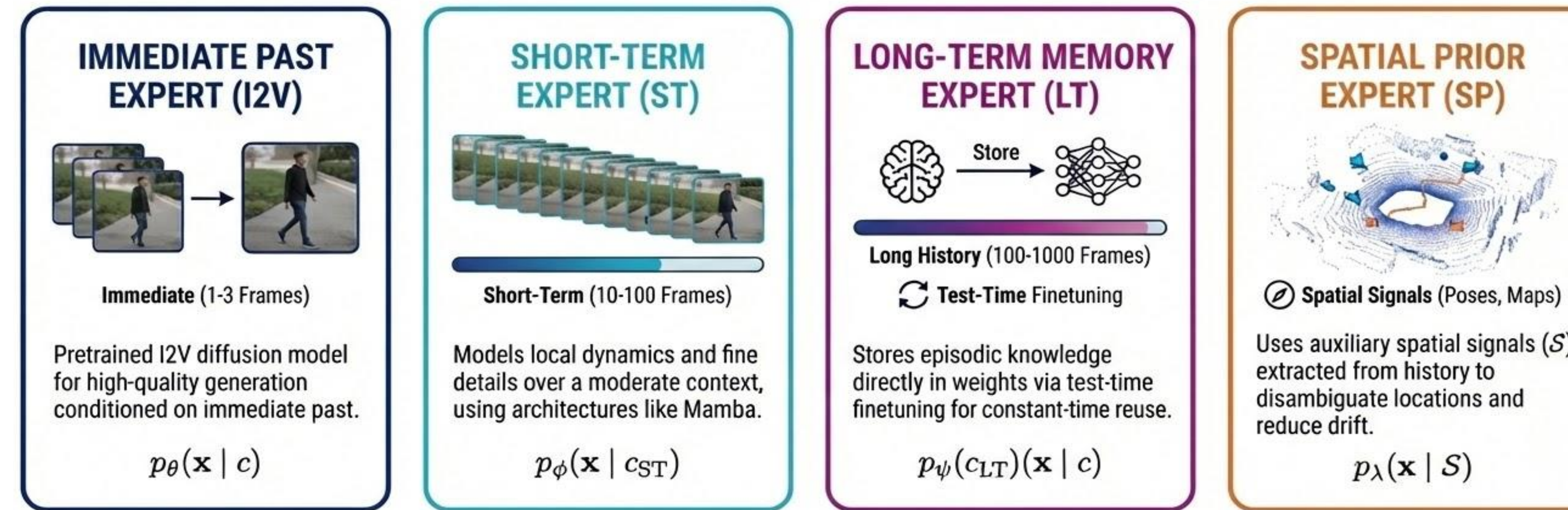


1. Building Consistent World Models

We aim to build world models that generate videos that are **consistent** with past observations in an **efficient and scalable** manner



4. Memory as Product of Experts



- Memory state composed of multiple contexts:
 - Short-term past / Long-term past
 - Possibly other structured memory formats
- Each context is handled by a **specialized diffusion model**

$$\Rightarrow p(\mathbf{x} | \mathcal{M}) = p(\mathbf{x} | c_1, \dots, c_K) \propto \prod_{i=1}^K p_i(\mathbf{x} | c_i)$$

2. Different Approaches to Memory

Transformer	Recurrent Models	TT-Training	Explicit Memory
Context window acts as memory	Information stored in some state vector	Information stored in temporary parameters	Observations stored in an external database
$\mathcal{M}_t = \{K_{1:t}, V_{1:t}\}$ $O_t = \text{softmax}(q_t K_{1:t}^\top) V_{1:t}$ $K_{t+1} = \text{concat}(K_{1:t}, k_{t+1})$ $V_{t+1} = \text{concat}(V_{1:t}, v_{t+1})$	$\mathcal{M}_t = h_t$ $O_t = C \mathcal{M}_t$ $\mathcal{M}_{t+1} = A \mathcal{M}_t + B x_t$	$\mathcal{M}_t = \theta_t$ $O_t = f(\theta_t, x_t)$ $\mathcal{M}_{t+1} = \theta_{t+1} - \eta \nabla_{\theta} \mathcal{L}$	$\mathcal{M}_t = \{(e_i, d_i)\}_{i=1}^N$ $r_t = \text{TopK}(\text{sim}(q_t, E))$ $O_t = f(x_t, d_{r_t})$ $\mathcal{M}_{t+1} = \begin{cases} \mathcal{M}_t \cup (e_{t+1}, d_{t+1}) \\ \mathcal{M}_t \end{cases}$
<ul style="list-style-type: none"> • ✓ Enables highly detailed, accurate predictions • ✗ Expensive: scales data, model size, and test-time compute 	<ul style="list-style-type: none"> • ✓ Cheap reading and writing into hidden states • ✗ Limited by state capacity 	<ul style="list-style-type: none"> • ✓ Higher capacity than fixed hidden states • ✗ More expensive update mechanism 	<ul style="list-style-type: none"> • ✓ Cheap reading and exact retrieval • ✗ Memory grows, search grows

5. Key Enabler

Combine independent diffusion models **at sampling time** by adding up the scores to sample from the product distribution

$$x_0 \sim p_{\theta}(x_0)$$

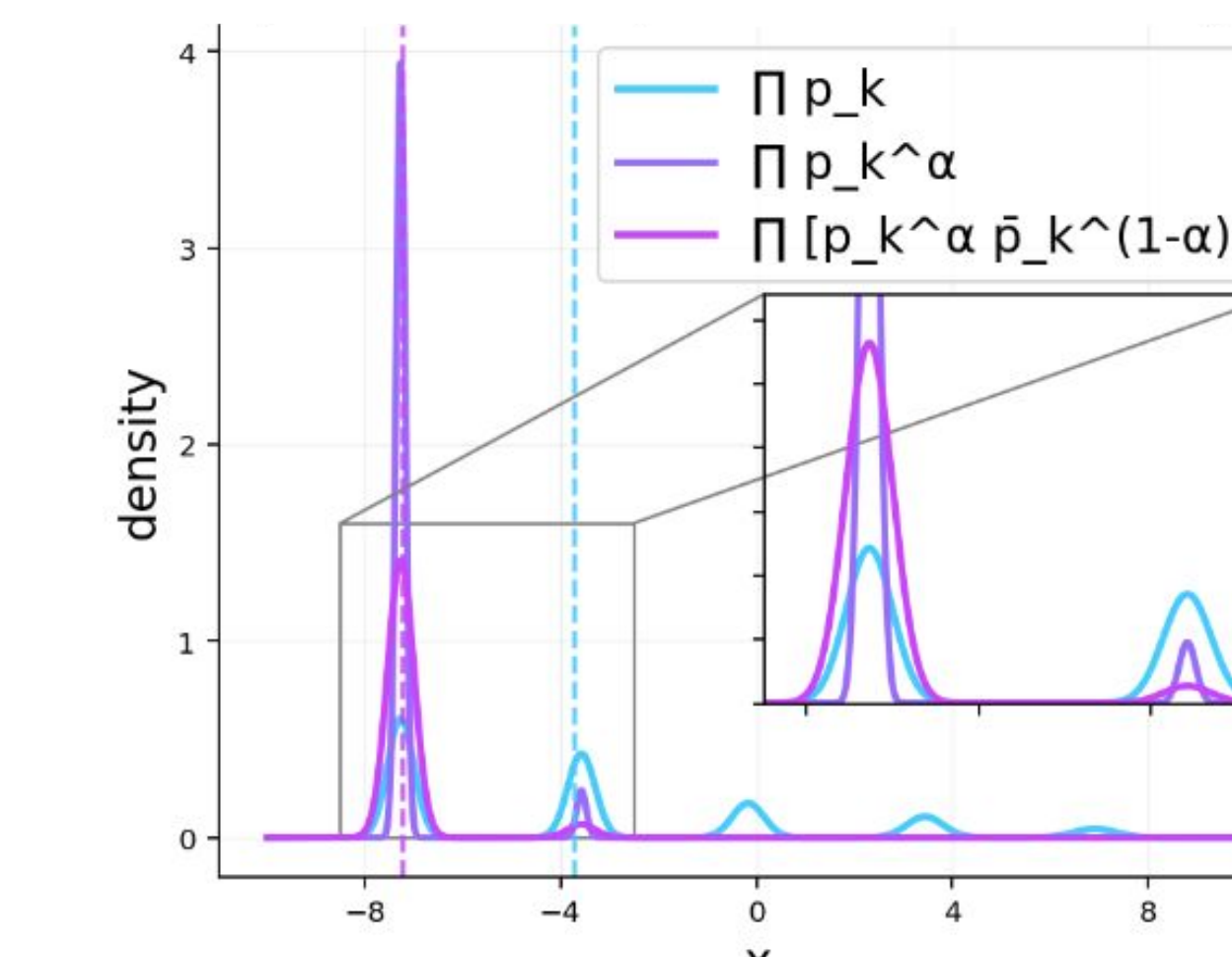
$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \nabla_{x_t} \log p_{\theta}(\mathbf{x}_t, t)) + \sigma_t \eta$$

$$x_0 \sim p_{\theta}(x_0) p_{\phi}(x_0) / Z$$

$$\Rightarrow \mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - [\nabla_{x_t} \log p_{\theta}(\mathbf{x}_t, t) + \nabla_{x_t} \log p_{\phi}(\mathbf{x}_t, t)])$$

6. Spurious Modes

Spurious likelihood modes get amplified if simultaneously present in different experts \Rightarrow We can remove them in the CFG way:

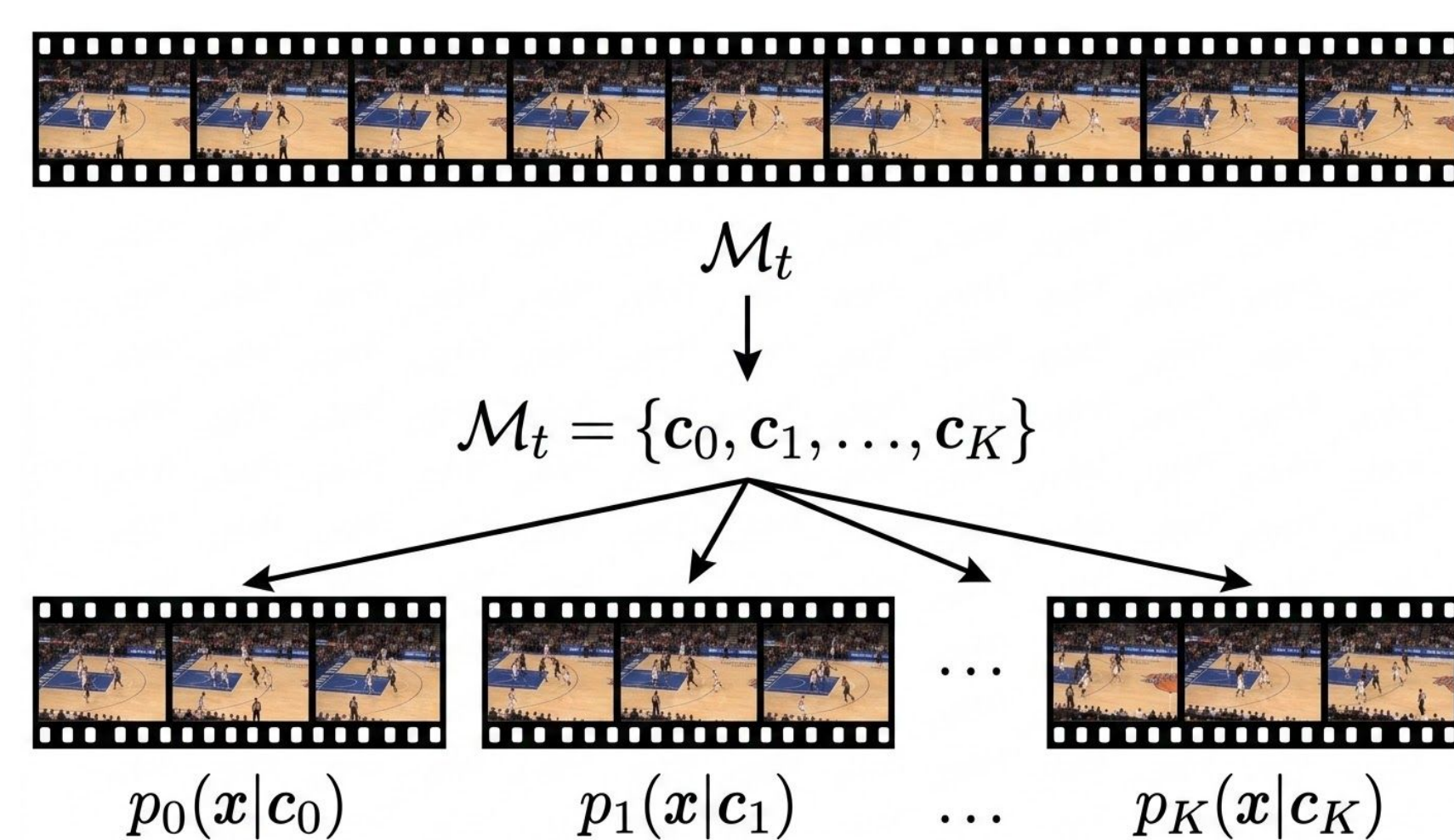


$$p_{\text{CoME}}(\mathbf{x} | c, c_{\text{ST}}, c_{\text{LT}}, \mathcal{S}) \propto [p_{\theta}(\mathbf{x} | \emptyset)^{1-\alpha_0} p_{\theta}(\mathbf{x} | c)^{\alpha_0}] \times [p_{\phi}(\mathbf{x} | \emptyset)^{1-\alpha_1} p_{\phi}(\mathbf{x} | c_{\text{ST}})^{\alpha_1}] \times [p_{\psi}(\emptyset)(\mathbf{x} | c)^{1-\alpha_2} p_{\psi}(c_{\text{LT}})(\mathbf{x} | c)^{\alpha_2}] \times [p_{\lambda}(\mathbf{x} | \emptyset)^{1-\alpha_3} p_{\lambda}(\mathbf{x} | \mathcal{S})^{\alpha_3}]$$

3. Can We Combine Them?

Memory could be distributed, not centralized

- Do **not** force one backbone to solve all temporal scales
- Use a **composition of specialized memory experts** for different parts of the history



Accurate predictions

Can the composition of experts make more accurate predictions?

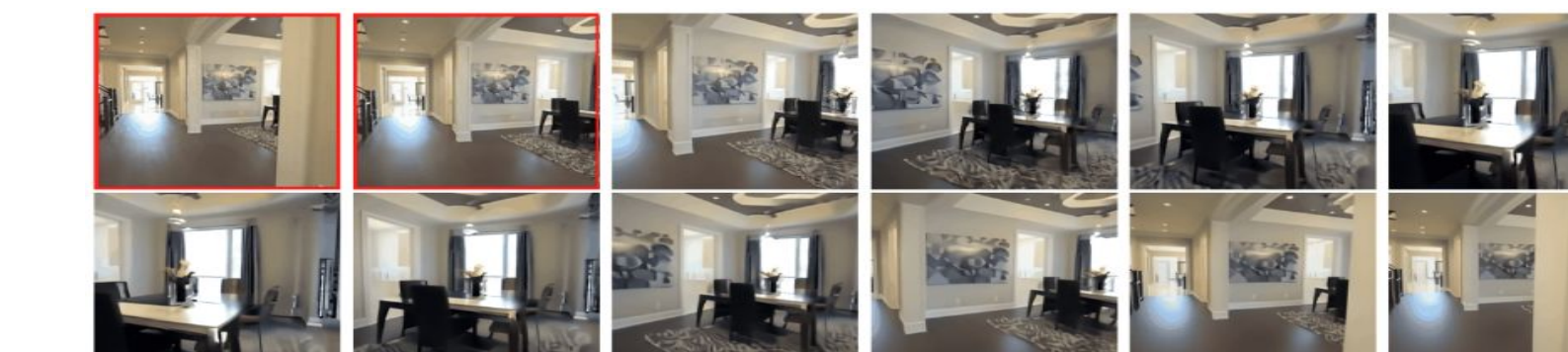
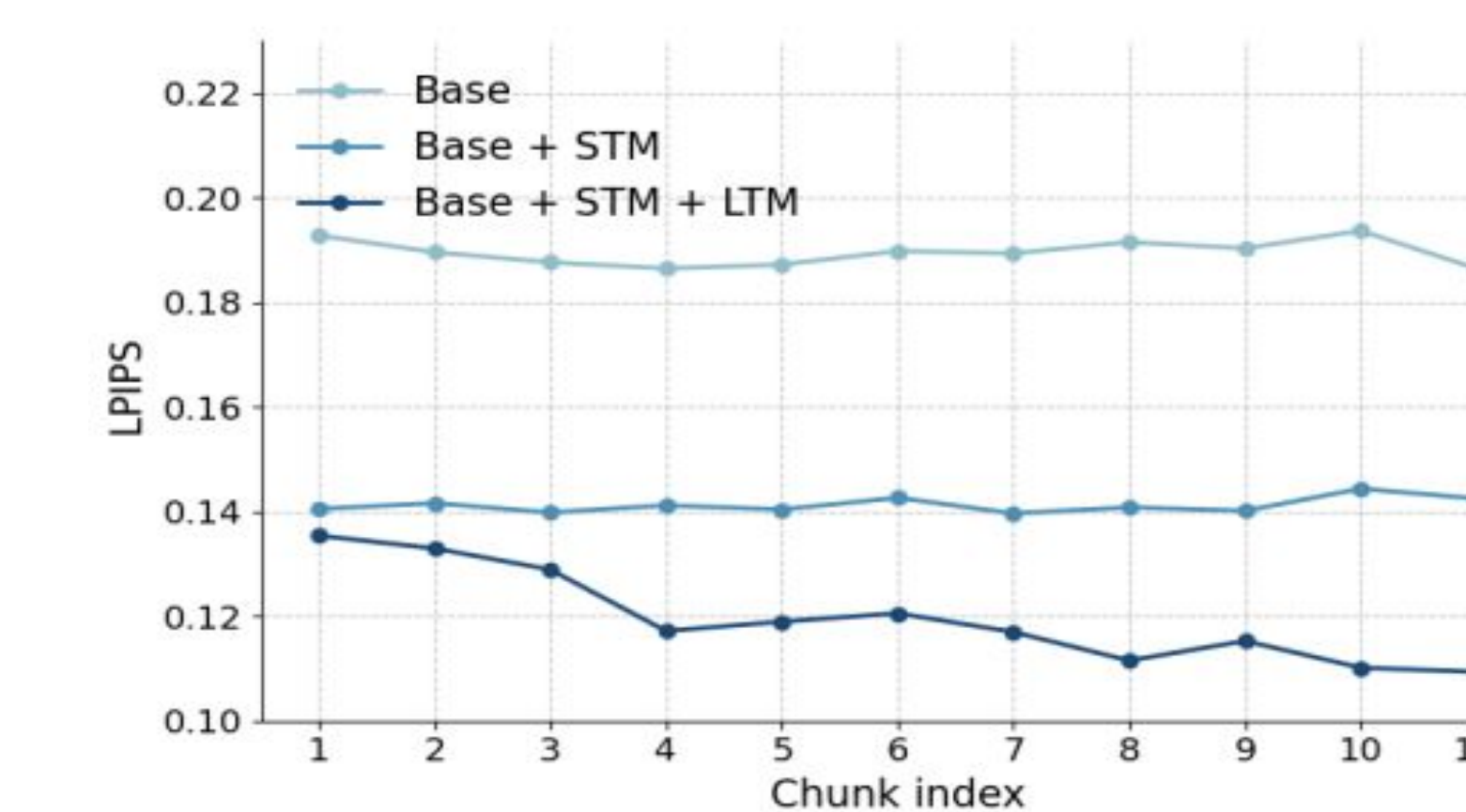
Method	LPIPS ↓	SSIM ↑	PSNR ↑
Base	0.209	0.771	19.16
+ STM	0.156	0.820	21.29
+ LTM	0.171	0.805	19.98
+ SLTM	0.150	0.833	20.65
+ STM+LTM	0.114	0.862	22.32
CoME	0.097	0.892	23.07
Sliding	0.183	0.753	19.02
SSM	0.158	0.828	20.62
Full	0.113	0.859	22.78



7. Experiments

Long-Term Memory

Does the TT-Learning of the long term memory contribute?



Planning Capability

Does the composition help with planning? We investigate the NWM setting:

	GNM	NOMAD	NWM	CoME
ATE (↓)	1.87	1.93	1.13	0.96
RPE (↓)	0.73	0.52	0.35	0.28

STM	LTM	NWM+STM	NWM+LTM
1.05	1.10	0.98	1.07
0.32	0.32	0.30	0.33

