



Contributions

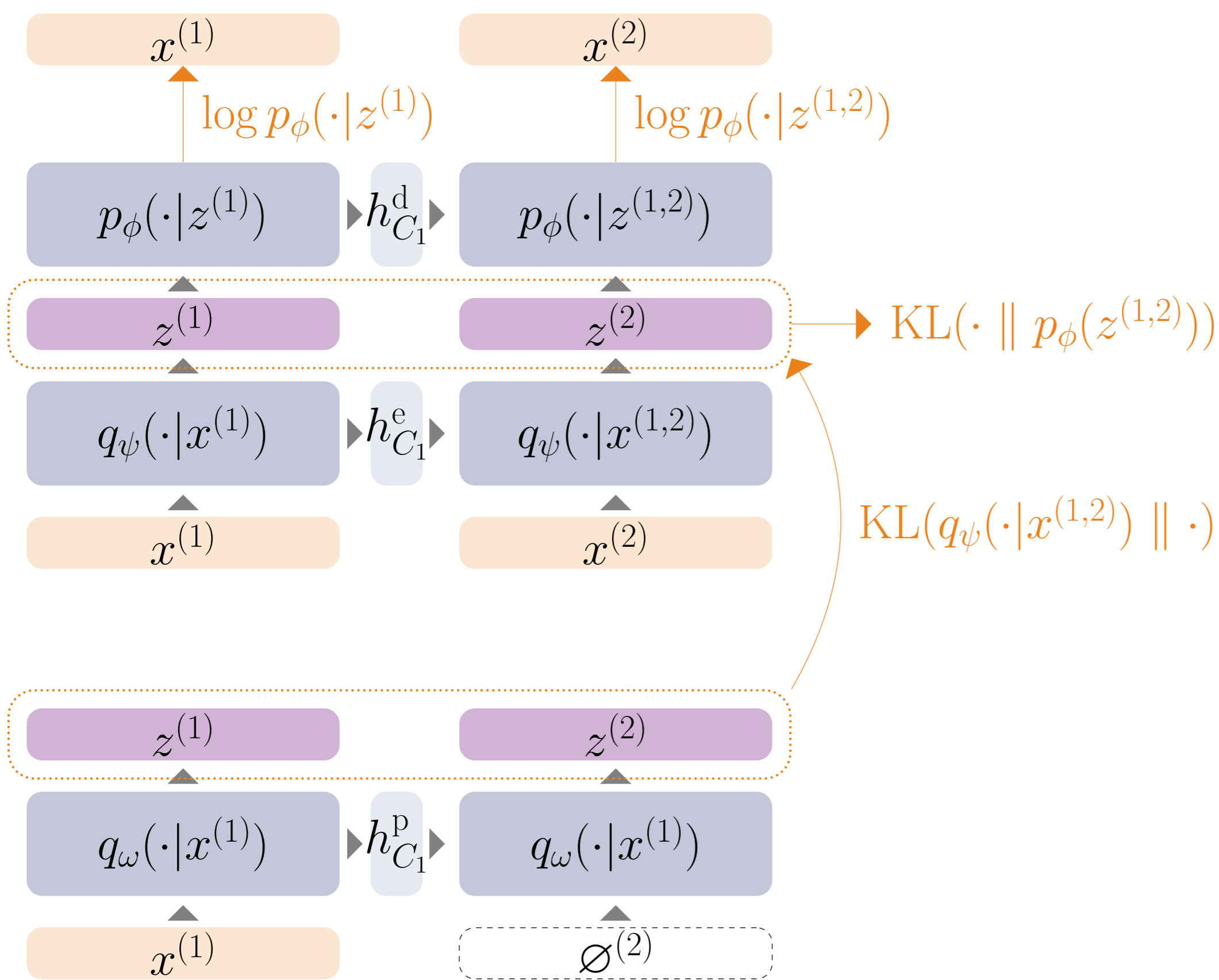
Summary

1. New VAE architecture with SSM encoder and decoder, enabling both **parallel training and generation**.
2. Introduction of a partial encoder based on an SSM allows **conditioning on partial realizations** (i.e., prompts).
3. Recurrence of the partial encoder and decoder allows **resuming generation without reprocessing** inputs.

Model	Generation	//	Prompt	Resume
Transformer	$O(T^2)$	✗	✓	✗
RNN	$O(T)$	✗	✓	✓
SSM	$O(T)$	✗	✓	✓
TVAE	$O(T^2)$	✓	✓	✗
VRNN	$O(T)$	✗	✓	✓
VSSM	$O(T)$	✓	✓	✓

Table 1. Model properties.

Training



1. Train posterior $q_\psi(z_{1:T}|x_{1:T})$ and generative $p_\phi(x_{1:T}|z_{1:T})$ according to the evidence lower bound (ELBO):

$$\max_{\phi, \psi} \mathbb{E}_{p(x_{1:T})} \left[\underbrace{\mathbb{E}_{q_\psi} \log p_\phi(x_{1:T}|z_{1:T}) - \text{KL}(q_\psi(z_{1:T}|x_{1:T}) \| p_\phi(z_{1:T}))}_{\text{ELBO}_{\phi, \psi}(x_{1:T})} \right].$$

2. Train partial posterior $q_\omega(z_{1:T}|x_{1:C})$, for $C \sim \mathcal{U}([0, T])$:

$$\min_{\omega} \mathbb{E}_{p(x_{1:C})} \text{KL}(p_\phi(z_{1:T}|x_{1:C}) \| q_\omega(z_{1:T}|x_{1:C})),$$

where,

$$\begin{aligned} p_\phi(z_{1:T}|x_{1:C}) &= \mathbb{E}_{p(x_{1:T}|x_{1:C})} p_\phi(z_{1:T}|x_{1:T}), \\ &\approx \mathbb{E}_{p(x_{1:T}|x_{1:C})} q_\psi(z_{1:T}|x_{1:T}). \end{aligned}$$

Generation

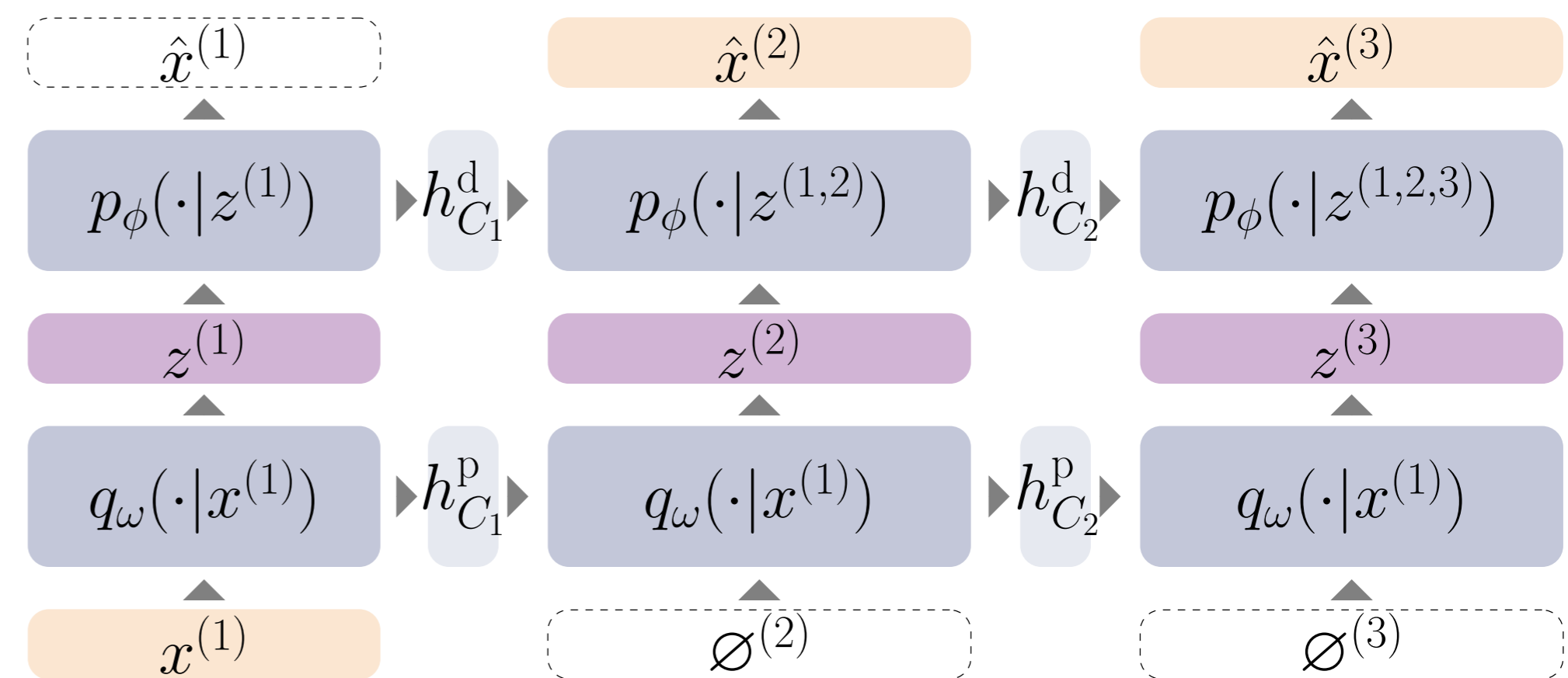


Fig 1. Given a prompt $x^{(1)} = x_{1:C_1}$, generate $\hat{x}^{(2)} = \hat{x}_{C_1:C_2}$ in parallel, then resume generation to produce $\hat{x}^{(3)} = \hat{x}_{C_2:C_3}$.

Architecture

Variational State Space Model (VSSM)

- **Prior distribution** uniform on a discrete latent space:

$$p_\phi(z_{1:T}) = \prod_{t=1}^T p_\phi(z_t) = \prod_{t=1}^T \frac{1}{N^Z},$$

- **Generative distribution** based on an SSM f_ϕ^{dec} :

$$p_\phi(x_{1:T}|z_{1:T}) = \prod_{t=1}^T \mathcal{P}(x_t | f_\phi^{\text{dec}}(z_{1:t})),$$

- **Posterior distribution** based on an SSM f_ψ^{enc} :

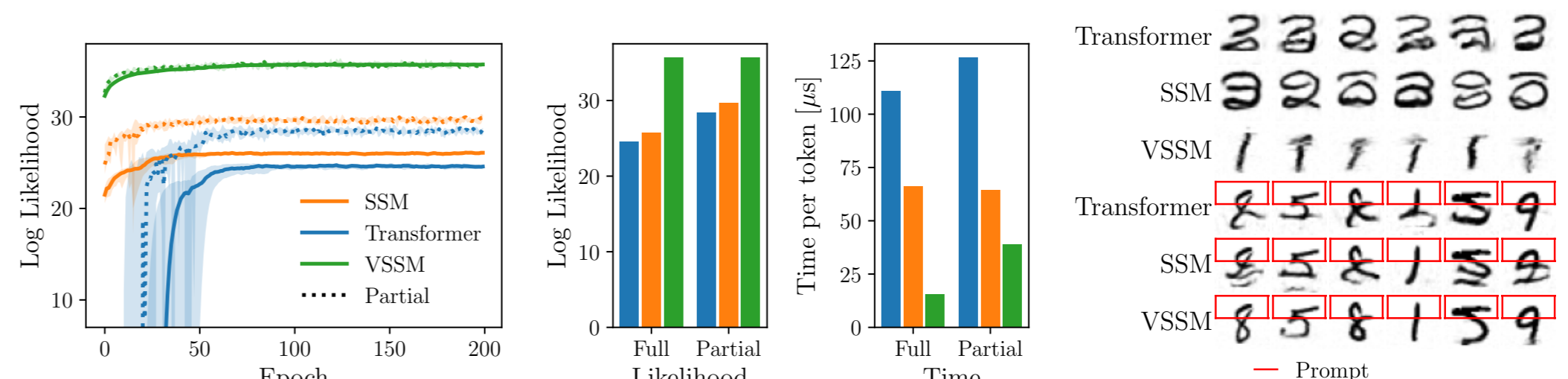
$$q_\psi(z_{1:T}|x_{1:T}) = \prod_{t=1}^T \mathcal{D}(z_t | f_\psi^{\text{enc}}(x_{1:t})),$$

- **Partial posterior distribution** based on an SSM f_ω^{par} :

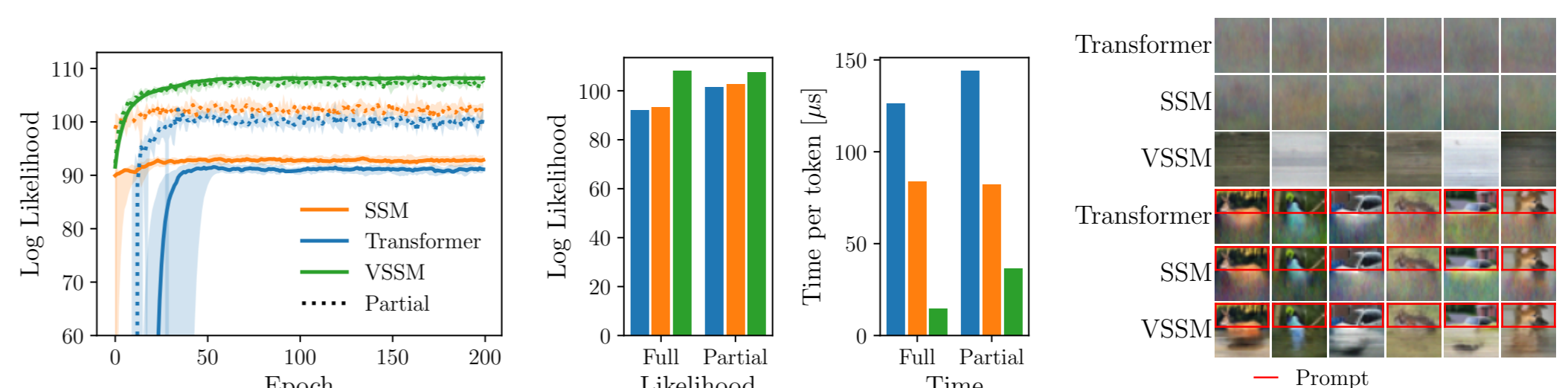
$$q_\omega(z_{1:T}|x_{1:C}) = \prod_{t=1}^T \mathcal{D}(z_t | f_\omega^{\text{par}}(\bar{x}_{1:t})),$$

where $\bar{x}_{1:T} = (x_{1:C}, \emptyset, \dots, \emptyset)$.

Results



(a) MNIST training, statistics and samples.



(b) CIFAR training, statistics and samples.

Fig 2. Validation likelihood (full and partial), test statistics (likelihood and time), samples (unconditional and prompted).