

Multi Time Scale World Models

Anonymous Author(s)
 Affiliation
 Address
 email

1 A Implementation Details

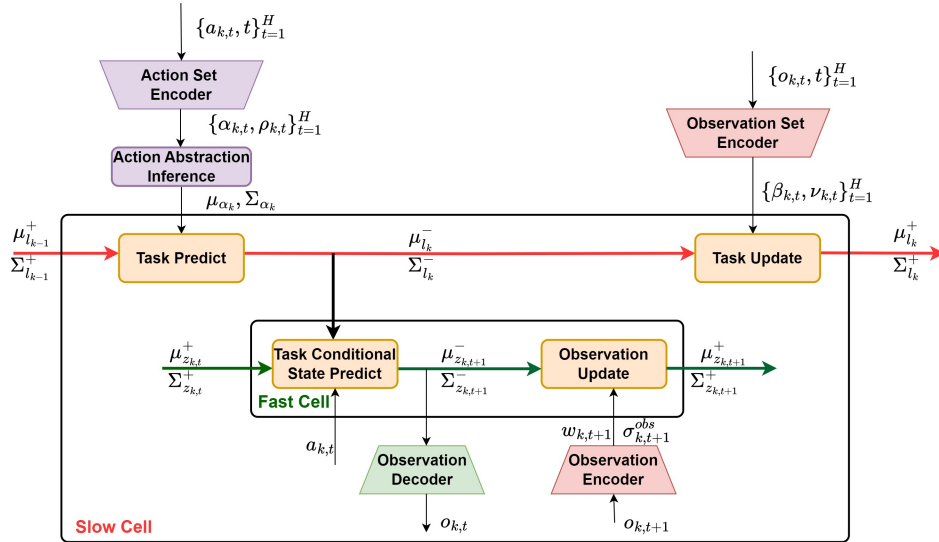


Figure 1: Schematic of a 2-Level MTS3 Architecture. Inference in MTS3 takes place via closed-form equations derived using exact inference, spread across two-time scales. For the fast time scale (fts) SSM, these include the task conditional state predict and observation update stages as discussed in Section 3.2 of the main paper. Whereas, for the slow time scale (sts) SSM, these include the task prediction and task update stages which are described in Section 3.3.

2 A.1 Inference In Slow Time Scale SSM

3 A.1.1 Inferring Action Abstraction (sts-SSM)

4 Given a set of encoded primitive actions and their correspond-
 5 ing variances $\{\alpha_{k,t}, \rho_{k,t}\}_{t=1}^H$, using the prior and observation
 6 model assumptions in Section 3.1.2 of main paper, we infer
 7 the latent abstract action $p(\alpha_k | \alpha_{k,1:H}) = \mathcal{N}(\mu_{\alpha_k}, \Sigma_{\alpha_k}) =$
 8 $\mathcal{N}(\mu_{\alpha_k}, \text{diag}(\sigma_{\alpha_k}))$ as a Bayesian aggregation [10] of these using
 9 the following closed-form equations:

$$\sigma_{\alpha_k} = \left((\sigma_0)^\ominus + \sum_{n=1}^N ((\rho_{k,t})^\ominus) \right)^\ominus,$$

$$\mu_{\alpha_k} = \mu_0 + \sigma_{\alpha_k} \odot \sum_{n=1}^N (\alpha_{k,t} - \mu_0) \odot \rho_{k,t}$$

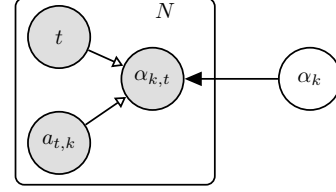


Figure 2: Generative model for the abstract action α_k . The hollow arrows are deterministic transformations leading to implicit distribution $\alpha_{k,t}$ using an action set encoder.

10 Here, \ominus , \odot and \otimes denote element-wise inversion, product, and
 11 division, respectively. The update equation is coded as the ‘‘abstract action inference’’ neural network
 12 layer as shown in Figure 1.

13 A.1.2 Task Prediction (sts-SSM)

14 The goal of this step is to update the prior marginal over the latent task variable l_k , $p(l_k | \beta_{1:k-1}, \alpha_{1:k})$,
 15 given the posterior beliefs from the time window $k - 1$ and abstract action α_k .

16 Using the linear dynamics model assumptions from Section 3.3, we can use the following closed-form
 17 update equations to compute, $p(l_k | \beta_{1:k-1}, \alpha_{1:k}) = \mathcal{N}(\mu_{l_k}^-, \Sigma_{l_k}^-)$, where

$$\begin{aligned} \mu_{l_k}^- &= \mathbf{X} \mu_{l_{k-1}}^+ + \mathbf{Y} \alpha_k \\ \Sigma_{l_k}^- &= \mathbf{X} \Sigma_{l_{k-1}}^+ \mathbf{X}^T + \mathbf{Y} \Sigma_{\alpha_k} \mathbf{Y}^T + \mathbf{S}. \end{aligned} \quad (1)$$

18 These closed-form equations are coded as the ‘‘task predict’’ neural net layer as shown in Figure 1.

19 A.1.3 Task Update (sts-SSM)

20 In this stage, we update the prior over l_k using an abstract observation set $\{\beta_{k,t}\}_{t=1}^H$, to ob-
 21 tain the latent task the posterior $\mathcal{N}(\mu_{z_{k,t}}^+, \Sigma_{z_{k,t}}^+) = \mathcal{N}\left(\begin{bmatrix} \mu_t^{u+} \\ \mu_t^{l+} \end{bmatrix}, \begin{bmatrix} \Sigma_t^u & \Sigma_t^s \\ \Sigma_t^s & \Sigma_t^l \end{bmatrix}^+\right)$, with $\Sigma_{l_k}^u =$
 22 $\text{diag}(\sigma_{l_k}^u)$, $\Sigma_{l_k}^l = \text{diag}(\sigma_{l_k}^l)$ and $\Sigma_{l_k}^s = \text{diag}(\sigma_{l_k}^s)$.

23 To do so we first invert the prior covariance matrix $\begin{bmatrix} \Sigma_{l_k}^u & \Sigma_{l_k}^s \\ \Sigma_{l_k}^s & \Sigma_{l_k}^l \end{bmatrix}^+$ to the precision matrix
 24 $\begin{bmatrix} \lambda_{l_k}^u & \lambda_{l_k}^s \\ \lambda_{l_k}^s & \lambda_{l_k}^l \end{bmatrix}^+$ for permutation invariant parallel processing. The posterior precision is then com-
 25 puted using scalar operations are follows, where only $\lambda_{l_k}^u$ is changed by

$$\lambda_{l_k}^{u+} = \lambda_{l_k}^{u-} + \sum_{t=1}^H \mathbf{1} \otimes \nu_{k,t} \quad (2)$$

26 while $\lambda_{l_k}^{l+} = \lambda_{l_k}^{l-}$ and $\lambda_{l_k}^{s+} = \lambda_{l_k}^{s-}$ remain constant. The operator \otimes denotes the element-wise
 27 division. The posterior precision is inverted back to the posterior covariance vectors $\sigma_{l_k}^{u+}$, $\sigma_{l_k}^{l+}$ and
 28 $\sigma_{l_k}^{s+}$. Now, the posterior mean $\mu_{l_k}^+$ can be obtained from the prior mean $\mu_{l_k}^-$ as

$$\mu_{l_k}^+ = \mu_{l_k}^- + \begin{bmatrix} \sigma_{l_k}^{u+} \\ \sigma_{l_k}^{s+} \end{bmatrix} \odot \left[\begin{array}{c} \sum_{t=1}^H (\beta_{k,t} - \mu_{l_k}^{u,-}) \otimes \nu_{k,t} \\ \sum_{t=1}^H (\beta_{k,t} - \mu_{l_k}^{s,-}) \otimes \nu_{k,t} \end{array} \right]. \quad (3)$$

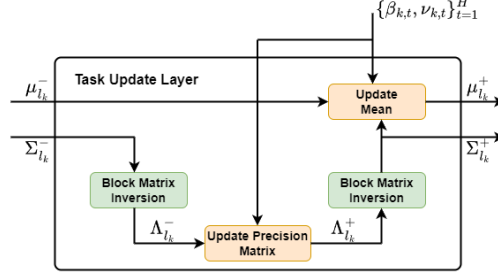


Figure 3: Implementation of task update layer which performs posterior latent task inference in the sts-SSM.

29

30 The inversion between the covariance matrix and precision matrix can be done via scalar operations
 31 leveraging block diagonal structure as derived in Appendix B. Figure 3 shows the schematic of the
 32 task update layer.

33 A.2 Inference In Fast Time Scale SSM

34 The inference in fts-SSM for a time-window k involves two stages as illustrated in Figure ??,
 35 calculating the prior and posterior over the latent state variable z_t . To keep the notation uncluttered,
 36 we will also omit the time-window index k whenever the context is clear as in section 3.2.

37 A.2.1 Task Conditional State Prediction (fts-SSM)

38 Following the assumptions of a task conditional linear dynamics as in Section 3.2 of the main paper,
 39 we obtain the prior marginal for $p(z_{k,t} | w_{1:t-1}^k, a_{1:t-1}^k, \beta_{1:k-1}, \alpha_{1:k-1}) = \mathcal{N}(\mu_{z_{k,t}}^-, \Sigma_{z_{k,t}}^-)$ in closed
 40 form, where

$$\begin{aligned} \mu_{z_{k,t}}^- &= \mathbf{A} \mu_{z_{k,t-1}}^- + \mathbf{B} a_{k,t-1} + \mathbf{C} \mu_{l_k}^-, \\ \Sigma_{z_{k,t}}^- &= \mathbf{A} \Sigma_{z_{k,t-1}}^+ \mathbf{A}^T + \mathbf{C} \Sigma_{l_k}^- \mathbf{C}^T + \mathbf{Q}. \end{aligned} \quad (4)$$

41 A.2.2 Observation Update (fts-SSM)

42 In this stage, we compute the posterior belief $p(z_{k,t} | w_{1:t}^k, a_{1:t}^k, \beta_{1:k}, \alpha_{1:k-1}) = \mathcal{N}(\mu_{z_{k,t}}^-, \Sigma_{z_{k,t}}^-)$.
 43 using the same closed-form update as in [1]. The choice of the special observation model
 44 splits the state into two parts, an upper z_t^u and a lower part z_t^l , resulting in the posterior belief

$$45 \mathcal{N}(\mu_{z_{k,t}}^-, \Sigma_{z_{k,t}}^-) = \mathcal{N}\left(\begin{bmatrix} \mu_t^{u+} \\ \mu_t^{l+} \end{bmatrix}, \begin{bmatrix} \Sigma_t^u & \Sigma_t^s \\ \Sigma_t^s & \Sigma_t^l \end{bmatrix}^+\right), \text{ with } \Sigma_t^u = \text{diag}(\sigma_t^s), \Sigma_t^l = \text{diag}(\sigma_t^l) \text{ and}$$

46 $\Sigma_t^s = \text{diag}(\sigma_t^s)$. Thus, the factorization allows for only the diagonal and one off-diagonal vector of
 47 the covariance to be computed and simplifies the calculation of the mean and posterior to simple
 48 scalar operations.

49 The closed-form equations for the mean can be expressed as the following scalar equations,

$$z_t^+ = z_t^- + \begin{bmatrix} \sigma_t^{u,-} \\ \sigma_t^{l,-} \end{bmatrix} \odot \begin{bmatrix} w_t - z_t^{u,-} \\ w_t - z_t^{l,-} \end{bmatrix} \oslash \begin{bmatrix} \sigma_t^{u,-} + \sigma_t^{\text{obs}} \\ \sigma_t^{u,-} + \sigma_t^{\text{obs}} \end{bmatrix},$$

50 The corresponding equations for the variance update can be expressed as the following scalar
 51 operations,

$$\begin{aligned} \sigma_t^{u,+} &= \sigma_t^{u,-} \odot \sigma_t^{u,-} \oslash (\sigma_t^{u,-} + \sigma_t^{\text{obs}}), \\ \sigma_t^{s,+} &= \sigma_t^{u,-} \odot \sigma_t^{s,-} \oslash (\sigma_t^{u,-} + \sigma_t^{\text{obs}}), \\ \sigma_t^{l,+} &= \sigma_t^{l,-} - \sigma_t^{s,-} \odot \sigma_t^{s,-} \oslash (\sigma_t^{u,-} + \sigma_t^{\text{obs}}), \end{aligned}$$

52 , where \odot denotes the elementwise vector product and \oslash denotes an elementwise vector division.

53 **A.3 Modelling Assumptions**

54 **A.3.1 Control Model**

55 To achieve action conditioning within the recurrent cell of fts-SMM, we include a control model
 56 $b(a_{k,t})$ in addition to the linear transition model A_t . $b(a_{k,t}) = f(a_{k,t})$, where $f(\cdot)$ can be any
 57 non-linear function approximator. We use a multi-layer neural network regressor with ReLU activa-
 58 tions [8].

59 However, unlike the fts-SSM where actions are assumed to be known and subjected to no noise, in
 60 the sts-SSM, the abstract action is an inferred latent variable with an associated uncertainty estimate.
 61 Hence we use a linear control model Y , for principled uncertainty propagation.

62 **A.3.2 Transition Noise**

63 We assume the covariance of the transition noise Q and S in both timescales to be diagonal. The
 64 noise is learned and is independent of the latent state.

65 **A.4 Training**

66 **A.4.1 Training Objective Derivation**

67 We further expand on the training objective in Section 4.2 here. The training objective for the MTS3
 68 involves maximizing the posterior predictive log-likelihood which for a single trajectory, can be
 69 derived as,

$$\begin{aligned}
 L &= \sum_{k=1}^N \sum_{t=1}^H \log p(\mathbf{o}_{k,t+1} | \boldsymbol{\beta}_{1:k-1}, \boldsymbol{\alpha}_{1:k-1}, \mathbf{w}_{k,1:t}, \mathbf{a}_{k,1:t}) \\
 &= \sum_{k=1}^N \sum_{t=1}^H \log \int \int p(\mathbf{o}_{k,t+1} | \mathbf{z}_{k,t+1}) p(\mathbf{z}_{k,t+1} | \mathbf{w}_{k,1:t}, \mathbf{a}_{k,1:t}, \mathbf{l}_k) p(\mathbf{l}_k | \boldsymbol{\beta}_{1:k-1}, \boldsymbol{\alpha}_{1:k-1}) d\mathbf{z}_{k,t+1} d\mathbf{l}_k \\
 &= \sum_{k=1}^N \sum_{t=1}^H \log \int p(\mathbf{o}_{k,t+1} | \mathbf{z}_{k,t+1}) p_{\mathbf{l}_k}(\mathbf{z}_{k,t+1} | \mathbf{w}_{k,1:t}, \mathbf{a}_{k,1:t}) d\mathbf{z}_{k,t+1}. \tag{5}
 \end{aligned}$$

70 The extension to multiple trajectories is straightforward. The approximation to the objective is done
 71 based on a moment-matching perspective as discussed in Section 4.2 of the main paper.

72 **A.4.2 Initialization**

73 We initialize the states \mathbf{l}_1 and $\mathbf{z}_{1,1}$ at both timescales for the first-time window $k = 1$ with an all zeros
 74 vector and corresponding covariance matrices as $\boldsymbol{\Sigma}_{\mathbf{l}_1} = \boldsymbol{\Sigma}_{\mathbf{z}_{1,1}} = 10 \cdot \mathbf{I}$. For subsequent windows, the
 75 prior belief $p(\mathbf{z}_{k,1})$ for the first time step of time window k , is initialized using the posterior belief
 76 $p_{\mathbf{l}_{k-1}}(\mathbf{z}_{k-1,H} | \mathbf{w}_{k-1,1:H}, \mathbf{a}_{k-1,1:H})$ of the last time step of time window $k - 1$.

77 It is also crucial to correctly initialize the transition matrix at both time scales so that the transition
 78 does not yield an unstable system. Initially, the transition model should focus on copying the encoder
 79 output so that the encoder can learn how to extract good features if observations are available and
 80 useful. We initialize the diagonal elements of the transition matrix at both timescales with 1 and the
 81 off-diagonal elements with 0.2, while the rest of the elements are set to 0, a choice inspired from [1].

82 **A.4.3 Learnable Parameters**

83 The learnable parameters in the computation graph are as follows:

84 **Fast Time Scale SSM:** The linear transition model A , the non-linear control factor b , the linear
 85 latent task transformation model C , the transition noise Q , along with the observation encoder and
 86 the output decoder.

87 **Slow Time Scale SSM:** The linear transition model X , the linear control model Y , the transition
 88 noise S , along with the observation set encoder and the action set encoder.

89 **B Proofs and Derivations**

90 In the following sections vectors are denoted by a lowercase letter in
 91 bold, such as " \mathbf{v} ", while Matrices as an uppercase letter in bold, such
 92 as " \mathbf{M} ". \mathbf{I} denotes identity matrix and $\mathbf{0}$ represents a matrix filled with
 93 zeros. For any matrix \mathbf{M} , \mathbf{m} denotes the corresponding vector of diagonal
 94 entries. Also, \odot denotes the elementwise vector product and \oslash denotes
 95 an elementwise vector division.

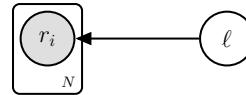


Figure 4: Graphical Model For Bayesian conditioning with N observations.

96 **B.1 Bayesian Conditioning As Permutation Invariant Set**
 97 **Operations**

98 **Gaussian Update Rule 1** (Bayesian Conditioning). *Consider the graphical model given in Figure*
 99 *4, where a set of N conditionally i.i.d observations $\bar{\mathbf{r}} = \{\mathbf{r}_i\}_{i=1}^N$ are generated by a latent variable*
 100 *\mathbf{l} and the observation model $p(\mathbf{r}_i|\mathbf{l}) = \mathcal{N}(\mathbf{r}_i | \mathbf{H}\mathbf{l}, \text{diag}(\sigma_i^{obs}))$. Assuming an observation model*
 101 *$\mathbf{H} = [\mathbf{I}, \mathbf{0}]$, the mean ($\boldsymbol{\mu}$) and precision matrix ($\boldsymbol{\Lambda}$) of the posterior over the latent variable \mathbf{l} , $p(\mathbf{l}|\bar{\mathbf{r}}) =$
 102 *$\mathcal{N}(\boldsymbol{\mu}_l^+, \boldsymbol{\Sigma}_l^+) = \mathcal{N}(\boldsymbol{\mu}_l^+, (\boldsymbol{\Lambda}_l^+)^{-1})$, given the prior $p_0(\mathbf{l}) = \mathcal{N}(\boldsymbol{\mu}_l^-, \boldsymbol{\Sigma}_l^-) = \mathcal{N}(\boldsymbol{\mu}_l^-, (\boldsymbol{\Lambda}_l^-)^{-1})$ have*
 103 *the following permutation invariant closed form updates.**

$$\begin{aligned} \boldsymbol{\Lambda}_l^+ &= \boldsymbol{\Lambda}_l^- + \begin{bmatrix} \text{diag}(\sum_{i=1}^N \frac{1}{\sigma_i^{obs}}), & \mathbf{0} \\ \mathbf{0}, & \mathbf{0} \end{bmatrix} \\ \boldsymbol{\mu}_l^+ &= \boldsymbol{\mu}_l^- + \begin{bmatrix} \boldsymbol{\sigma}_l^{u+} \\ \boldsymbol{\sigma}_l^{s+} \end{bmatrix} \odot \begin{bmatrix} \sum_{i=1}^N (\mathbf{r}_i - \boldsymbol{\mu}_l^{u,-}) \odot \frac{1}{\sigma_i^{obs}} \\ \sum_{i=1}^N (\mathbf{r}_i - \boldsymbol{\mu}_l^{s,-}) \odot \frac{1}{\sigma_i^{obs}} \end{bmatrix} \end{aligned} \quad (6)$$

Note that $\boldsymbol{\Sigma}_l$ is the covariance matrix which is the inverse of the precision matrix $\boldsymbol{\Lambda}_l$. Due to the observation model assumption $\mathbf{H} = [\mathbf{I}, \mathbf{0}]$, they take block diagonal form,

$$\boldsymbol{\Sigma}_l = \begin{bmatrix} \boldsymbol{\Sigma}_l^u & \boldsymbol{\Sigma}_l^s \\ \boldsymbol{\Sigma}_l^s & \boldsymbol{\Sigma}_l^l \end{bmatrix}, \text{ with } \boldsymbol{\Sigma}_l^u = \text{diag}(\boldsymbol{\sigma}_l^u), \boldsymbol{\Sigma}_l^l = \text{diag}(\boldsymbol{\sigma}_l^l) \text{ and } \boldsymbol{\Sigma}_l^s = \text{diag}(\boldsymbol{\sigma}_l^s).$$

104 **Proof:**

105 **Case 1 (Single Observation):** Before deriving the update rule for N conditionally iid observations,
 106 let us start with a simpler case consisting of a single observation \mathbf{r} . If the marginal Gaussian distri-
 107 bution for the latent variable \mathbf{l} takes the form $p(\mathbf{l}) = \mathcal{N}(\mathbf{l} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$ and the conditional Gaussian
 108 distribution for the single observation \mathbf{r} given \mathbf{l} has the form, $p(\mathbf{r} | \mathbf{l}) = \mathcal{N}(\mathbf{r} | \mathbf{H}\mathbf{l} + \mathbf{b}, \mathbf{L}^{-1})$. Then
 109 the posterior distribution over \mathbf{l} can be obtained in closed form as,

$$p(\mathbf{l} | \mathbf{r}) = \mathcal{N}(\mathbf{l} | \boldsymbol{\Sigma} \{ \mathbf{H}^T \mathbf{L} \mathbf{r} + \boldsymbol{\Lambda} \boldsymbol{\mu} \}, \boldsymbol{\Lambda}^{-1}), \text{ where } \boldsymbol{\Lambda} = (\boldsymbol{\Lambda} + \mathbf{H}^T \mathbf{L} \mathbf{H}). \quad (7)$$

110 We refer to [2] to the proof for this standard result.

Case 2 (Set Of Observations): Now instead of a single observation, we wish to derive a closed form
 solution for the posterior over latent variable $\mathbf{l} \in \mathbb{R}^{2d}$, given a set of N conditionally i.i.d observations
 $\bar{\mathbf{r}} = \{\mathbf{r}_i\}_{i=1}^N$. Here each element $\mathbf{r}_i \in \mathbb{R}^d$ of the set $\bar{\mathbf{r}}$ is assumed to to have an observation model
 $\mathbf{H} = [\mathbf{I}, \mathbf{0}]$. In the derivation, we represent the set of N observations as a random vector

$$\bar{\mathbf{r}} = \begin{bmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \\ \cdot \\ \cdot \\ \mathbf{r}_N \end{bmatrix}_{Nd \times 1}.$$

Since each observation in the set $\bar{\mathbf{r}}$ are conditionally independent, we denote the conditional distri-
 bution over the context set as $\bar{\mathbf{r}} | \mathbf{l} \sim \mathcal{N}(\bar{\mathbf{H}}\mathbf{l}, \boldsymbol{\Sigma}_r)$, where the diagonal covariance matrix has the
 following form:

$$\boldsymbol{\Sigma}_r = \begin{bmatrix} \text{diag}(\boldsymbol{\sigma}_{r_1}), & 0, & 0, & \dots, & 0 \\ 0, & \text{diag}(\boldsymbol{\sigma}_{r_2}), & 0, & \dots, & 0 \\ \cdot, & \cdot, & \cdot, & \cdot, & \cdot \\ \cdot, & \cdot, & \cdot, & \cdot, & \cdot \\ 0, & 0, & 0, & \dots, & \text{diag}(\boldsymbol{\sigma}_{r_N}) \end{bmatrix}_{Nd \times Nd}.$$

The corresponding observation model $\bar{\mathbf{H}}$ is

$$\bar{\mathbf{H}} = \begin{bmatrix} \mathbf{H} \\ \mathbf{H} \\ \cdot \\ \cdot \\ \mathbf{H} \end{bmatrix}_{Nd \times 2d} = \begin{bmatrix} \mathbf{I}, \mathbf{0} \\ \mathbf{I}, \mathbf{0} \\ \cdot, \cdot \\ \cdot, \cdot \\ \mathbf{I}, \mathbf{0} \end{bmatrix}_{Nd \times 2d}.$$

111 Now given the prior over the latent task variable $l \sim \mathcal{N}(\mu_l^-, \Sigma_l^-)$, the parameters of the posterior
 112 distribution over the task variable, $p(l|\bar{r}) \sim \mathcal{N}(\mu_l^+, \Lambda_l^+)$, can be obtained in closed-form substituting
 113 in Equation (7) as follows.

$$\begin{aligned} \Lambda_l^+ &= (\Sigma_l^+)^{-1} \\ &= \Sigma_l^{-1} + \bar{\mathbf{H}}^T \Sigma_r \bar{\mathbf{H}} \\ &= \Sigma_l^{-1} + \begin{bmatrix} \text{diag}(\sigma_{r_1}), \text{diag}(\sigma_{r_2}), \text{diag}(\sigma_{r_3}), \dots, \text{diag}(\sigma_{r_N}) \\ \mathbf{0}, \quad \mathbf{0}, \quad \mathbf{0}, \quad \dots, \quad \mathbf{0} \end{bmatrix}_{2d \times nd} \bar{\mathbf{H}} \\ &= \lambda_l^- + \begin{bmatrix} \text{diag}(\sum_{i=1}^n \frac{1}{\sigma_{r_i}}), \quad \mathbf{0} \\ \mathbf{0}, \quad \mathbf{0} \end{bmatrix}_{2d \times 2d} \end{aligned}$$

$$\begin{aligned} \mu_l^+ &= \mu_l^- + (\Lambda^+)^{-1} \bar{\mathbf{H}}^T (\sigma_r^{-2} \mathbf{I}) (\mathbf{y} - \bar{\mathbf{H}} \mu_x) \\ &= \mu_l^- + \Sigma^+ \bar{\mathbf{H}} (\sigma_r^{-2} \mathbf{I}) (\mathbf{y} - \bar{\mathbf{H}} \mu_x) \\ &= \mu_l^- + \Sigma^+ \begin{bmatrix} \sigma_{r_1}^{-2} \mathbf{I}, \sigma_{r_2}^{-2} \mathbf{I}, \sigma_{r_3}^{-2} \mathbf{I}, \dots, \sigma_{r_n}^{-2} \mathbf{I} \\ \mathbf{0}, \quad \mathbf{0}, \quad \mathbf{0}, \quad \dots, \quad \mathbf{0} \end{bmatrix} (\mathbf{y} - \bar{\mathbf{H}} \mu_x) \\ &= \mu_l^- + \begin{bmatrix} \sigma_l^{u+}, & \sigma_l^{s+} \\ \sigma_l^{s+}, & \sigma_l^{l+} \end{bmatrix} \begin{bmatrix} \sum_{n=1}^N (\mathbf{r}_n - \mu_l^{u,-}) \odot \frac{1}{\sigma_i} \\ \mathbf{0} \end{bmatrix} \\ &= \mu_l^- + \begin{bmatrix} \sigma_l^{u+} \\ \sigma_l^{s+} \end{bmatrix} \odot \begin{bmatrix} \sum_{i=1}^N (\mathbf{r}_i - \mu_l^{u,-}) \odot \frac{1}{\sigma_{r_i}} \\ \sum_{i=1}^N (\mathbf{r}_n - \mu_l^{u,-}) \odot \frac{1}{\sigma_{r_i}} \end{bmatrix} \end{aligned} \tag{8}$$

114 Here μ_l^+ is the posterior mean and Λ_l^+ is the posterior precision matrix.

115 **Corollary 1.** *The closed form updates for the resulting posterior distribution $p(l|\bar{r})$ is permutation*
 116 *invariant with respect to the observation set \bar{r} .*

117 B.2 Derivation For Matrix Inversions as Scalar Operations

118 **Inversion Of Block Diagonal Matrix.** *Consider a block matrix of the following form $\mathbf{A} =$*
 119 $\begin{bmatrix} \text{diag}(\mathbf{a}^u) & \text{diag}(\mathbf{a}^s) \\ \text{diag}(\mathbf{a}^s) & \text{diag}(\mathbf{a}^l) \end{bmatrix}$. *Then inverse $\mathbf{A}^{-1} = \mathbf{B}$ can be calculated using scalar operations*

120 *and is given as, $\mathbf{B} = \begin{bmatrix} \text{diag}(\mathbf{b}^u) & \text{diag}(\mathbf{b}^s) \\ \text{diag}(\mathbf{b}^s) & \text{diag}(\mathbf{b}^l) \end{bmatrix}$ where,*

$$\begin{aligned} \mathbf{b}^u &= \mathbf{a}^l \oslash (\mathbf{a}^u \odot \mathbf{a}^l - \mathbf{a}^s \odot \mathbf{a}^s) \\ \mathbf{b}^s &= -\mathbf{a}^s \oslash (\mathbf{a}^u \odot \mathbf{a}^l - \mathbf{a}^s \odot \mathbf{a}^s) \\ \mathbf{b}^l &= \mathbf{a}^u \oslash (\mathbf{a}^u \odot \mathbf{a}^l - \mathbf{a}^s \odot \mathbf{a}^s) \end{aligned} \tag{9}$$

121 .

122 **Proof:** To prove this we will use the following matrix identity of a partitioned matrix from [2],
 123 which states

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix} \tag{10}$$

where \mathbf{M} is defined as

$$\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}.$$

124 Here \mathbf{M} is called the Schur complement of the Matrix on the left side of Equation 10. The algebraic
 125 manipulations to arrive at scalar operations in Equation 9 are straightforward.

126 C Metrics Used For Measuring Long Horizon Predictions

127 C.1 Sliding Window RMSE

The sliding window RMSE (Root Mean Squared Error) metric is computed for a predicted trajectory in comparison to its ground truth. At each time step, the RMSE for each trajectory is determined by taking the root mean square of the differences between the ground truth and predicted values within a sliding window that terminates at the current time step. This sliding window, with a specified size, provides a smoothed localized assessment of prediction accuracy over the entire prediction length. Mathematically, the sliding window RMSE at time step t is given by:

$$\text{RMSE}(t) = \sqrt{\frac{1}{W} \sum_{i=t-W+1}^t (\text{gt}_i - \text{pred}_i)^2}$$

128 where t is the current time step, W is the window size, and gt_i and pred_i are the ground truth and
129 predicted values at time step i , respectively. The extension to multiple trajectories is straightforward
130 and omitted to keep the notation uncluttered.

131 C.2 Sliding Window NLL

132 The sliding window NLL (Negative Log-Likelihood) metric is computed for a predicted probability
133 distribution against the true distribution. At each time step, the NLL is determined by summing
134 the negative log-likelihood values within a sliding window that terminates at the current time step.
135 This sliding window, with a specified size, provides a smoothed localized evaluation of prediction
136 accuracy across the entire sequence.

Mathematically, the sliding window NLL at time step t is given by:

$$\text{NLL}(t) = -\frac{1}{W} \sum_{i=t-W+1}^t \log \mathcal{N}(\text{gt}_i | \text{predMean}_i, \text{predVar}_i)$$

137 where t is the current time step, W is the window size. predMean_i , predVar_i , and gt_i represent the
138 predicted mean, predicted variance, and the ground truth at time step i .

139 D Additional Experiments and Plots

140 D.1 Additional results on ablation with discretization step $H \cdot \Delta t$

141 In addition to the Hydraulics Dataset discussed in Section
142 6.4, we report the results of the ablation study with different
143 values of $H \cdot \Delta t$, for the mobile robot dataset. The higher the
144 value of H , the slower the timescale of the task dynamics
145 relative to the state dynamics. As seen in Figure 5, smaller
146 values of H (like 2,3,5 and 10) give significantly worse
147 performance. Very large values of H (like 150) also result in
148 degradation of performance. In the paper, we used a value
149 of $H=75$.

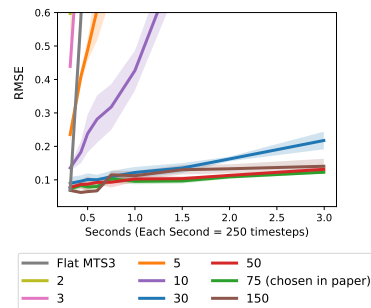


Figure 5: Ablation on discretization step $H \cdot \Delta t$. The long-term prediction results in terms of RMSE, with different H on the mobile dataset.

150 D.2 Visualization 151 of predictions given by different models.

152 In this section, we plot the multistep ahead predictions (mean
153 and variance) by different models on 3 datasets on normalized test trajectories. Not that we omit NaN
154 values in predictions while plotting.

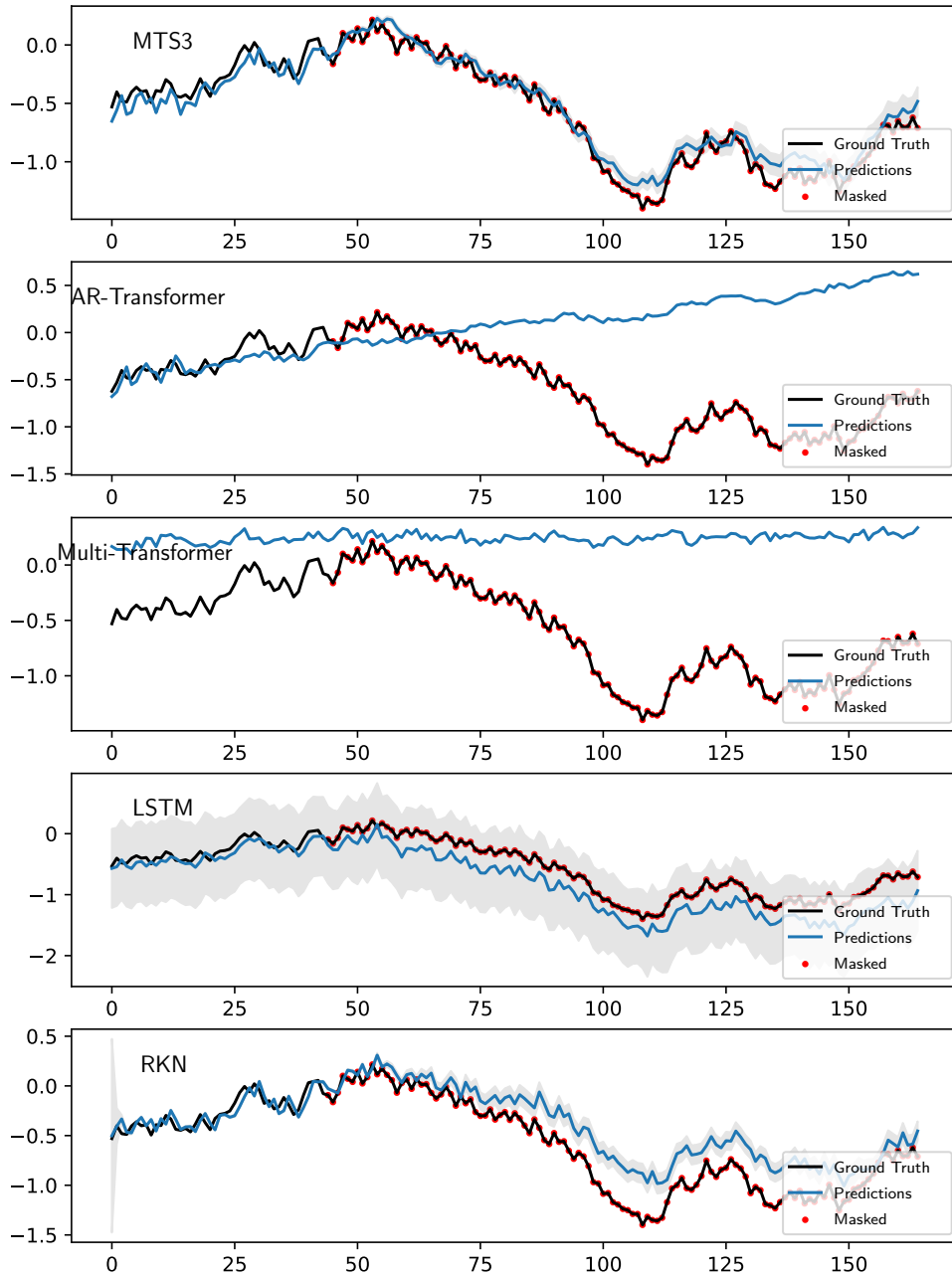


Figure 6: Multi-step ahead mean and variance predictions for a particular joint (joint 1) of Franka Kitchen Environment. The multi-step ahead prediction starts from the first red dot, which indicates masked observations. MTS3 gives the most reliable mean and variance estimates.

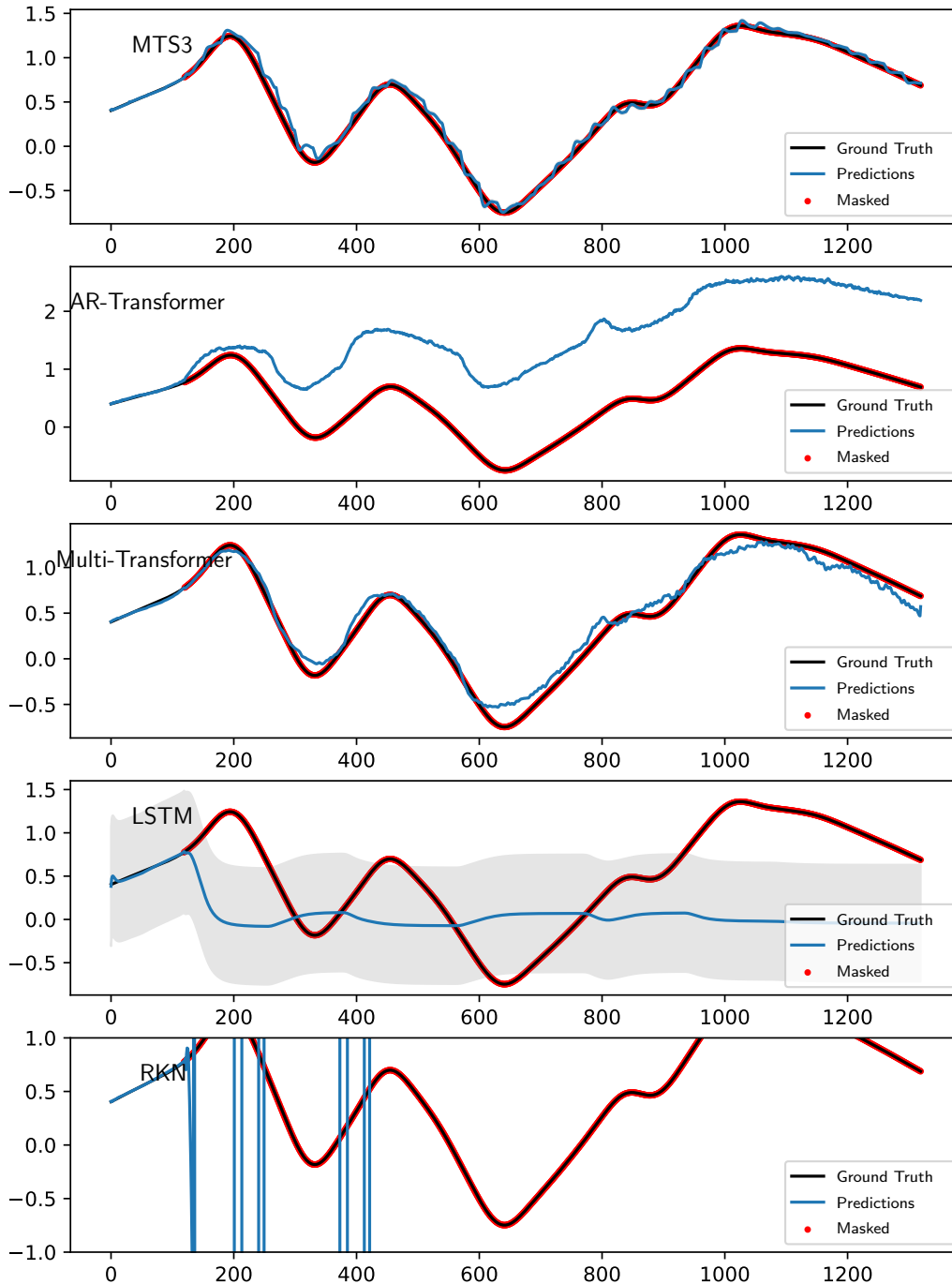


Figure 7: Multi-step ahead mean and variance predictions for a particular joint (joint 1) of Excavator Dataset. The multi-step ahead prediction starts from the first red dot, which indicates masked observations. MTS3 gives the most reliable mean and variance estimates even up to 12 seconds into the future. Another interesting observation can also be seen in the predictions for MTS3, where after every window k of sts-SSM, which is 0.3 seconds (30 timesteps) long, the updation of the higher-level abstractions helps in grounding the lower-level predictions thus helping in the long horizon yet fine-grained predictions.

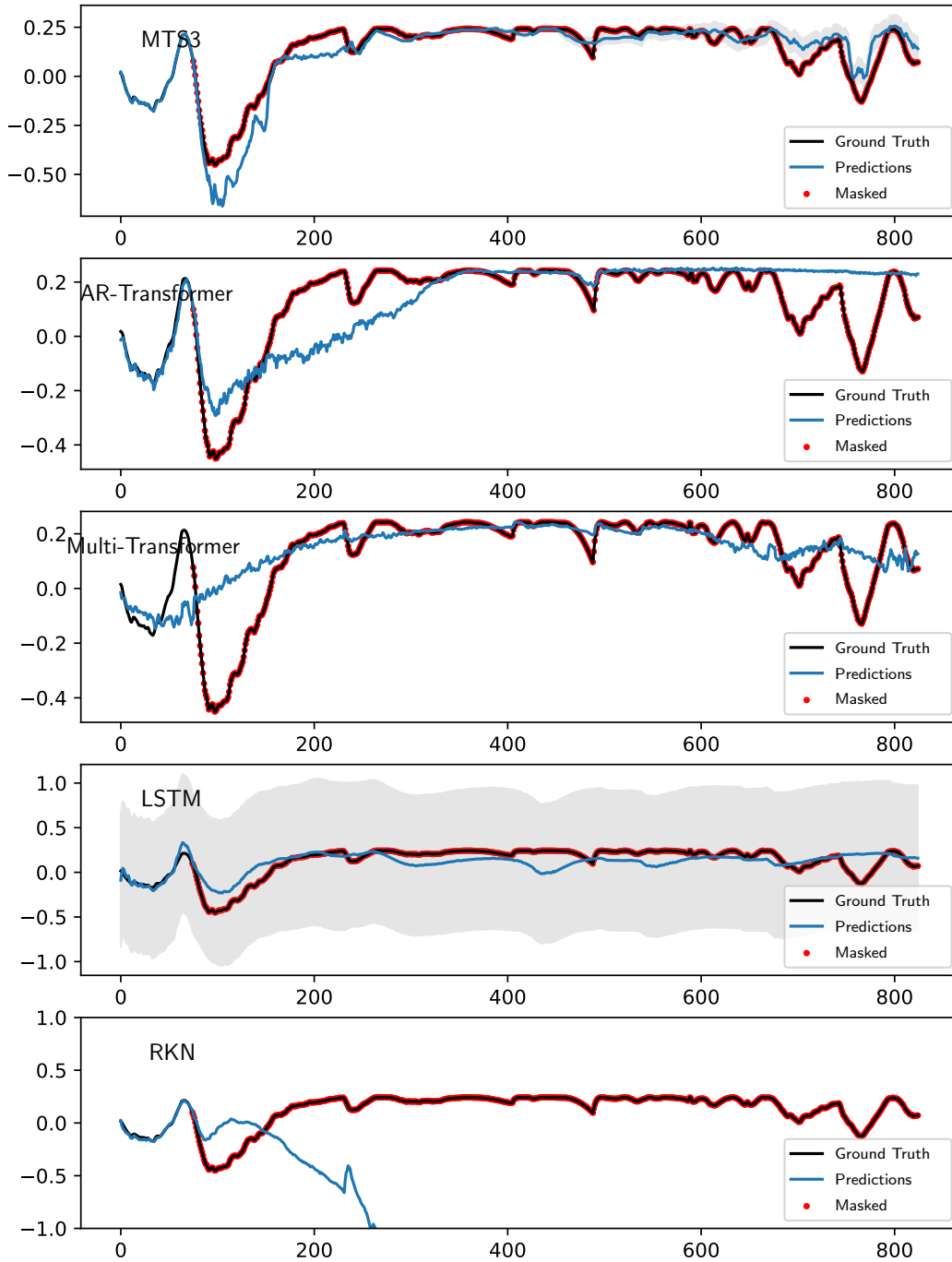


Figure 8: Multi-step ahead mean and variance predictions for a particular joint (joint 7) of Mobile Robot Dataset. The multi-step ahead prediction starts from the first red dot, which indicates masked observations. MTS3 gives the most accurate mean and variance estimates among all algorithms.

158 E Robots and Data

159 In all datasets, we only use information about agent/object positions and we mask out velocities
160 to create a partially observable setting. All datasets are subjected to a mean zero, unit variance
161 normalization during training. During testing, they are denormalized after predictions. The details of
162 the different datasets used are explained below:

163 E.1 D4RL Datasets

164 **Details:** We use a set of 3 different environments/agents from D4RL dataset [4], which includes
165 the HalfCheetah, Franka Kitchen and Maze2D (medium) environment. **(a) HalfCheetah:** We used
166 1000 suboptimal trajectories collected from a policy trained to approximately 1/3 the performance
167 of the expert. The observation space consists of 8 joint positions and the action space consists of 6
168 joint torques collected at 50 Hz frequency. 800 trajectories were used for training and 200 for testing.
169 For the long horizon task, we used 1.2 seconds (60 timesteps) as context and tasked the model to
170 predict 6 seconds (300 timesteps) into the future. **(b) Franka Kitchen:** The goal of the Franka
171 Kitchen environment is to interact with the various objects to reach a desired state configuration. The
172 objects you can interact with include the position of the kettle, flipping the light switch, opening and
173 closing the microwave and cabinet doors, or sliding the other cabinet door. We used the "complete"
174 version of the dataset and collected 1000 trajectories where all four tasks are performed in order. The
175 observation space consists of 30 dimensions (9 joint positions of the robot and 21 object positions).
176 The action space consists of 9 joint velocities clipped between -1 and 1 rad/s. The data was collected
177 at a 50 Hz frequency. 800 trajectories were used for training and 200 for testing. For the long horizon
178 task, we used 0.6 seconds (30 timesteps) as context and tasked the model to predict 2.7 seconds (135
179 timesteps) into the future. The dataset is complex due to multi-task, multi-object interactions in a
180 single trajectory. **(c) Medium Maze:** We used 20000 trajectories from a 2D Maze environment,
181 where each trajectory consists of a force-actuated ball (along the X and Y axis) moving to a fixed
182 target location. The observation consists of as the (x, y) locations and a 2D action space. The data is
183 collected at 100 Hz frequency. 16000 trajectories were used for training and 4000 for testing. For the
184 long horizon task, we used 0.6 seconds (60 timesteps) as context and tasked the model to predict 3.9
185 seconds (390 timesteps) into the future. Rendering of the three environments is shown in Figure 9.

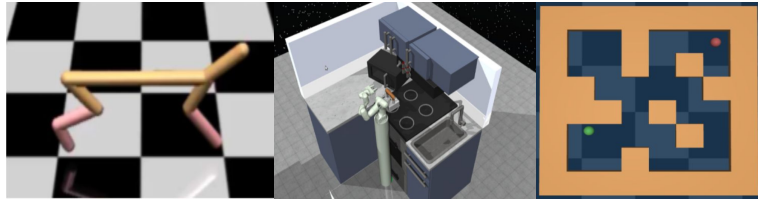


Figure 9: D4RL Environments: (left) HalfCheetah (middle) Franka Kitchen (right) Maze2D-Medium

186 E.2 Hydraulic Excavator

187 **Details:** We collected the data from a wheeled excavator JCB Hydradig 110W show in Figure
188 10. The data was collected by actuating the boom and arm of the excavator using Multisine and
189 Amplitude-Modulated Pseudo-Random Binary Sequence (APRBS) joystick signals with safety
190 mechanisms in place. A total of 150 mins of data was collected at a frequency of 100 Hz. of which
191 was used as a training dataset and the rest as testing. The observation space consists of the boom and
192 arm positions, while the joystick signals are chosen as actions. For the long horizon task we used 1.5
193 seconds (150 timesteps) as context and tasked the model to predict 12 seconds (1200 timesteps) into
194 the future.

195

196 E.3 Panda Robot With Varying Payloads

197 **Details:** We collected the data from a 7 DoF Franka Emika Panda manipulator during free motion
198 and while manipulating loads with weights 0kg (free motion), 0.5 kg, 1 kg, 1.5 kg, 2 kg and 2.5

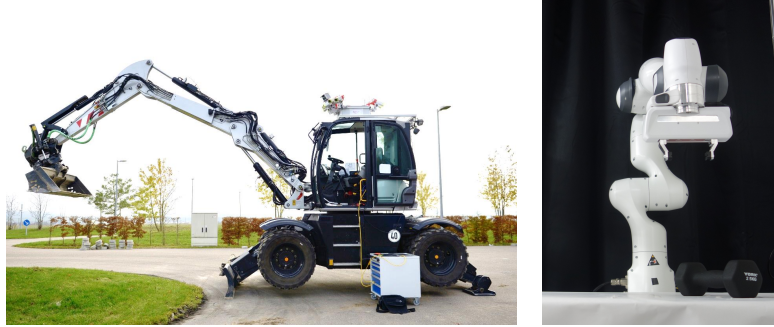


Figure 10: (left) JCB Hydradig 110W Excavator (right) Franka Emika Panda Robot

199 kg. The robot used is shown in Figure 10. Data is sampled at a frequency of 100 Hz. The training
 200 trajectories were motions with loads of 0kg (free motion), 1kg, 1.5kg, and 2.5 kgs, while the testing
 201 trajectories contained motions with loads of 0.5 kg and 2 kg. The observations for the forward model
 202 consist of the seven joint angles in radians, and the corresponding actions were joint Torques in Nm.
 203 For the long horizon task we used 0.6 seconds (60 timesteps) as context and tasked the model to
 204 predict 1.8 seconds (180 timesteps) into the future.

205

206 E.4 Wheeled Mobile Robot

Observation and Data Set: We collected 50 random trajectories from a Pybullet simulator a wheeled mobile robot traversing terrain with slopes generated by a mix of sine waves as shown in Figure 11. Data is sampled at high frequencies (500Hz). 40 out of the 50 trajectories were used for training and the rest 10 for testing. The observations consist of parameters which completely describe the location and orientation of the robot. The observation of the robot at any time instance t consists of the following features:

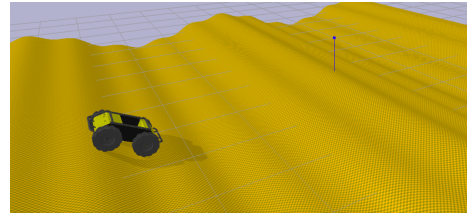


Figure 11: Wheeled Mobile Robot traversing terrain with complex variations in slopes induced by a mix of sine functions.

$$o_t = [x, y, z, \cos(\alpha), \sin(\alpha), \cos(\beta), \sin(\beta), \cos(\gamma), \sin(\gamma)]$$

207 where, x, y, z - denote the global position of the Center of Mass of the robot, α, β, γ - Roll, pitch
 208 and yaw angles of the robot respectively, in the global frame of reference [9]. For the long horizon
 209 task we used 0.6 seconds (150 timesteps) as context and tasked the model to predict 3 seconds (750
 210 timesteps) into the future.

211

212 **F Hyperparameters and Compute Resources**

213 **Compute Resources** For training MTS3, LSTM, GRU and Transformer models we used compute
214 nodes with (i) Nvidia 3090 and (ii) Nvidia 2080 RTX GPUs. For training more computationally
215 expensive locally linear models like RKN, HiP-RSSM we used compute nodes with NVIDIA A100-40
216 GPUs.

217 **Hyperparameters** Hyperparameters were selected via grid search. In general, the performance of
218 MTS3 is not very sensitive to hyperparameters. Among all the baselines, Transformer models were
219 most sensitive to hyperparameters (see Appendix E.5 for details of Transformer architecture).

220 **Discretization Step:** For MTS3, the discretization step for the slow time scale SSM as discussed in
221 Section 3.1 for all datasets was fixed as $\mathbf{H} \cdot \Delta t = 0.3$ seconds. In our experiments, we found that
222 discretization values between $0.2 \leq \mathbf{H} \cdot \Delta t \leq 0.5$ seconds give similar performance.

223 **Rule Of thumb for choosing discretization step in MTS3:** For any N-level MTS3 as defined in
224 Section 3.4, we recommend searching for discretization factor H_i as a hyperparameter. However,
225 as a general rule of thumb, it can be chosen as $H_i = (\sqrt[N]{T})^i$, where T is the maximum prediction
226 horizon required / episode length. This ensures that very long recurrences are divided between
227 smaller equal-length task-reconfigurable local SSM windows (of length $\sqrt[N]{T}$) spread across several
228 hierarchies.

229 **Encoder Decoder Architecture:** For all recurrent models (MTS3, HiP-RSSM, RKN, LSTM and
230 GRU) we use a similar encoder-decoder architecture across datasets. Small variations from these
231 encoder-decoder architecture hyperparameters can still lead to similar prediction performance as
232 reported in the paper.

233

234 Observation Set Encoder (MTS3): 1 fully connected + linear output:

- 235 • Fully Connected 1: 240, ReLU

236 Action Set Encoder (MTS3): 1 fully connected + linear output:

- 237 • Fully Connected 1: 240, ReLU

238 Observation Encoder (MTS3, HiP-RSSM, RKN, LSTM, GRU): 1 fully connected + linear output:

- 239 • Fully Connected 1: 120, ReLU

240 Observation Decoder (MTS3, HiP-RSSM, RKN, LSTM, GRU): 1 fully connected + linear output:

- 241 • Fully Connected 1: 120, ReLU

242 Control Model (Primitive Action Encoder) (MTS3, HiP-RSSM, RKN): 1 fully connected + linear
243 output:

- 244 • Fully Connected 1: 120, ReLU

245 The rest of the hyperparameters are described below:

246 **F.1 D4RL Datasets**

247 **F.1.1 Half Cheetah**

248 **Recurrent Models**

249 Transition Model (HiP-RSSM, RKN): number of basis: 32

- 250 • $\alpha(\mathbf{z}_t)$: No hidden layers - softmax output

Hyperparameters	MTS3	HiP-RSSM	RKN	LSTM	GRU
Learning Rate	3e-3	1e-3	1e-3	1e-3	1e-3
Latent Observation Dimension	15	15	15	15	15
Observation Set Latent Dimension (sts-SSM)	15	-	-	-	-
Latent State Dimension	30	30	30	45	45
Latent Task Dimension	30	30	-	-	-
Latent Abstract Action Dimension (sts-SSM)	30	-	-	-	-

251 **Autoregressive Transformer Baseline**

252 Learning Rate: 1e-5

253 Optimizer Used: Adam Optimizer

254 Embedding size: 96

255 Number of Decoder Layers: 4

256 Number Of Attention Heads: 4

257 **Multistep Transformer Baseline**

258 Learning Rate: 1e-5

259 Optimizer Used: Adam Optimizer

260 Embedding size: 128

261 Number Of Encoder Layers: 2

262 Number of Decoder Layers: 1

263 Number Of Attention Heads: 4

264 **F.1.2 Franka Kitchen**

Recurrent Models

Hyperparameters	MTS3	HiP-RSSM	RKN	LSTM	GRU
Learning Rate	3e-3	9e-4	9e-4	1e-3	1e-3
Latent Observation Dimension	30	30	30	30	30
Observation Set Latent Dimension (sts-SSM)	30	-	-	-	-
Latent State Dimension	60	60	60	90	90
Latent Task Dimension	60	60	-	-	-
Latent Abstract Action Dimension (sts-SSM)	60	-	-	-	-

265

266 Transition Model (HiP-RSSM, RKN): number of basis: 15

267 • $\alpha(z_t)$: No hidden layers - softmax output

268 **Autoregressive Transformer Baseline**

269 Learning Rate: 5e-5

270 Optimizer Used: Adam Optimizer

271 Embedding size: 64

272 Number of Decoder Layers: 4

273 Number Of Attention Heads: 4

274 **Multistep Transformer Baseline**

275 Learning Rate: 1e-5

276 Optimizer Used: Adam Optimizer

277 Embedding size: 64

278 Number Of Encoder Layers: 2

279 Number of Decoder Layers: 1

280 Number Of Attention Heads: 4

281 **F.1.3 Maze 2D**

282 **Recurrent Models**

283 Transition Model (HiP-RSSM, RKN): number of basis: 15

Hyperparameters	MTS3	HiP-RSSM	RKN	LSTM	GRU
Learning Rate	3e-3	9e-4	9e-4	1e-3	1e-3
Latent Observation Dimension	30	30	30	30	30
Observation Set Latent Dimension (sts-SSM)	30	-	-	-	-
Latent State Dimension	60	60	60	90	90
Latent Task Dimension	60	60	-	-	-
Latent Abstract Action Dimension (sts-SSM)	60	-	-	-	-

284 • $\alpha(z_t)$: No hidden layers - softmax output

285 **Autoregressive Transformer Baseline**

286 Learning Rate: 5e-5

287 Optimizer Used: Adam Optimizer

288 Embedding size: 96

289 Number of Decoder Layers: 4

290 Number Of Attention Heads: 4

291 **Multistep Transformer Baseline**

292 Learning Rate: 1e-5

293 Optimizer Used: Adam Optimizer

294 Embedding size: 128

295 Number Of Encoder Layers: 2

296 Number of Decoder Layers: 1

297 Number Of Attention Heads: 4

298 **F.2 Franka Robot Arm With Varying Loads**

Recurrent Models

Hyperparameters	MTS3	HiP-RSSM	RKN	LSTM	GRU
Learning Rate	3e-3	9e-4	9e-4	3e-3	3e-3
Latent Observation Dimension	15	15	15	15	15
Observation Set Latent Dimension (sts-SSM)	15	-	-	-	-
Latent State Dimension	30	30	30	45	45
Latent Task Dimension	30	30	-	-	-
Latent Abstract Action Dimension (sts-SSM)	30	-	-	-	-

299

300 Transition Model (HiP-RSSM,RKN): number of basis: 32

301 • $\alpha(z_t)$: No hidden layers - softmax output

302 **Autoregressive Transformer Baseline**

303 Learning Rate: 5e-5

304 Optimizer Used: Adam Optimizer

305 Embedding size: 64

306 Number of Decoder Layers: 4

307 Number Of Attention Heads: 4

308 **Multistep Transformer Baseline**

309 Learning Rate: 2e-5

310 Optimizer Used: Adam Optimizer

311 Embedding size: 64

312 Number Of Encoder Layers: 2

313 Number of Decoder Layers: 1

314 Number Of Attention Heads: 4

315 **F.3 Hydraulic Excavator**

316 Transition Model (HiP-RSSM,RKN): number of basis: 15

Hyperparameters	MTS3	HiP-RSSM	RKN	LSTM	GRU
Learning Rate	3e-3	8e-4	8e-4	1e-3	1e-3
Latent Observation Dimension	15	15	15	15	15
Observation Set Latent Dimension (sts-SSM)	15	-	-	-	-
Latent State Dimension	30	30	30	45	45
Latent Task Dimension	30	30	-	-	-
Latent Abstract Action Dimension (sts-SSM)	30	-	-	-	-

317 • coefficient net $\alpha(z_t)$: No hidden layers - softmax output

318 **Autoregressive Transformer Baseline**

319 Learning Rate: 1e-5
320 Optimizer Used: Adam Optimizer
321 Embedding size: 96
322 Number of Decoder Layers: 4
323 Number Of Attention Heads: 4

324 **Multistep Transformer Baseline**

325 Learning Rate: 5e-5
326 Optimizer Used: Adam Optimizer
327 Embedding size: 64
328 Number Of Encoder Layers: 2
329 Number of Decoder Layers: 1
330 Number Of Attention Heads: 4

331 **F.4 Wheeled Robot Traversing Uneven Terrain**

Hyperparameters	MTS3	HiP-RSSM	RKN	LSTM	GRU
Learning Rate	3e-3	8e-4	8e-4	1e-3	1e-3
Latent Observation Dimension	30	30	30	30	30
Observation Set Latent Dimension (sts-SSM)	30	-	-	-	-
Latent State Dimension	60	60	60	90	90
Latent Task Dimension	60	60	-	-	-
Latent Abstract Action Dimension (sts-SSM)	60	-	-	-	-

332 Transition Model (HiP-RSSM,RKN): number of basis: 15

333 • coefficient net $\alpha(z_t)$: No hidden layers - softmax output

334 **Autoregressive Transformer Baseline**

335 Learning Rate: 5e-5
336 Optimizer Used: Adam Optimizer
337 Embedding size: 128
338 Number of Decoder Layers: 4
339 Number Of Attention Heads: 4

340 **Multistep Transformer Baseline**

341 Learning Rate: 5e-5
342 Optimizer Used: Adam Optimizer
343 Embedding size: 64
344 Number Of Encoder Layers: 4
345 Number of Decoder Layers: 2
346 Number Of Attention Heads: 4

347 **F.5 Transformer Architecture Details**

348 For the AR-Transformer Baseline, we use a GPT-like autoregressive version of transformers except
349 that for the autoregressive input we also concatenate the actions to make action conditional predictions.

350 For Multi-Transformer we use the same direct multistep prediction and loss as in recent Transformer
 351 time-series forecasting literature [12, 6, 7, 11]. A description of the action conditional direct multi-step
 352 version of the transformer is given in Algorithm 1.

Algorithm 1: MultiStep Transformer

Require: Input past observations $\mathbf{o}_{inp} \in \mathbb{R}^{S \times C}$; Input Past Actions $\mathbf{a}_{inp} \in \mathbb{R}^{S \times A}$; Future
 Actions $\mathbf{a}_{pred} \in \mathbb{R}^{O \times A}$; Input Length S ; Predict length O ; Observation Dimension C ; Action
 Dimension A ; Feature dimension d_k ; Encoder layers number N ; Decoder layers number M .

```

1:  $\mathbf{o}_{inp} \in \mathbb{R}^{S \times C}$ ,  $\mathbf{a}_{inp} \in \mathbb{R}^{S \times A}$ ,  $\mathbf{a}_{pred} \in \mathbb{R}^{O \times A}$ 
2:  $\mathbf{X}_{inp} = \text{ConCatFeatureWise}(\mathbf{o}_{inp}, \mathbf{a}_{inp})$   $\triangleright \mathbf{X}_{inp} \in \mathbb{R}^{S \times (C+A)}$ 
3:  $\mathbf{X}_{pred} = \text{ConCatFeatureWise}(\text{Zeros}(O, C), \mathbf{a}_{pred})$   $\triangleright \mathbf{X}_{pred} \in \mathbb{R}^{O \times (C+A)}$ 
4:  $\mathbf{X}_{enc}, \mathbf{X}_{dec} = \text{ConCat}(\mathbf{X}_{inp}, \mathbf{X}_{pred})$   $\triangleright \mathbf{X}_{enc} \in \mathbb{R}^{S \times (C+A)}, \mathbf{X}_{dec} \in$   

 $\mathbb{R}^{(S+O) \times (C+A)}$ 
5:  $\mathbf{X}_{enc}^0 = \text{Embed}(\mathbf{X}_{enc})$   $\triangleright \mathbf{X}_{enc}^0 \in \mathbb{R}^{S \times d_k}$ 
6: for  $l$  in  $\{1, \dots, N\}$  do
  | 7:  $\mathbf{X}_{enc}^{l-1} = \text{LayerNorm}(\mathbf{X}_{enc}^{l-1} + \text{Attn}(\mathbf{X}_{enc}^{l-1}))$   $\triangleright \mathbf{X}_{enc}^{l-1} \in \mathbb{R}^{S \times d_k}$ 
  | 8:  $\mathbf{X}_{enc}^l = \text{LayerNorm}(\mathbf{X}_{enc}^{l-1} + \text{FFN}(\mathbf{X}_{enc}^{l-1}))$   $\triangleright \mathbf{X}_{enc}^l \in \mathbb{R}^{S \times d_k}$ 
end
9:  $\mathbf{X}_{dec}^0 = \text{Embed}(\mathbf{X}_{dec})$   $\triangleright \mathbf{X}_{dec}^0 \in \mathbb{R}^{(S+O) \times d_k}$ 
10: for for  $l$  in  $\{1, \dots, M\}$  do
  | 11:  $\mathbf{X}_{dec}^{l-1} = \text{LayerNorm}(\mathbf{X}_{dec}^{l-1} + \text{Attn}(\mathbf{X}_{dec}^{l-1}))$   $\triangleright$  Decoder
  | 12:  $\mathbf{X}_{dec}^{l-1} = \text{LayerNorm}(\mathbf{X}_{dec}^{l-1} + \text{Attn}(\mathbf{X}_{dec}^{l-1}, \mathbf{X}_{enc}^N))$   $\triangleright \mathbf{X}_{dec}^{l-1} \in \mathbb{R}^{(S+O) \times d_k}$ 
  | 13:  $\mathbf{X}_{dec}^l = \text{LayerNorm}(\mathbf{X}_{dec}^{l-1} + \text{FFN}(\mathbf{X}_{dec}^{l-1}))$   $\triangleright \mathbf{X}_{dec}^l \in \mathbb{R}^{(S+O) \times d_k}$ 
end
14:  $\mathbf{y} = \text{MLP}(\mathbf{X}_{dec}^M)$   $\triangleright \mathbf{y} \in \mathbb{R}^{(S+O) \times C}$ 
15: Return  $\mathbf{y}$   $\triangleright$  Return the prediction results

```

353 **G Limitations**

354 We list some of the limitations of the paper here. (i) We restricted our definition and experiments to
355 MTS3 with two levels of temporal abstractions, which was sufficient in many of our tasks. However,
356 for certain tasks like the Maze2D, we believe more hierarchies can help. As discussed in the main
357 paper the method and inference scheme allows easy addition of more Feudal [3] hierarchies with
358 larger discretization steps ($\mathbf{H} \cdot \Delta t$). (ii) We restrict our application to action conditional long horizon
359 future predictions and do not use the model for (hierarchical) control. A probabilistically principled
360 formalism for hierarchical control as an inference problem, that builds upon MTS3 models is left
361 for future work. (iii) Finally, we restrict our experiments to proprioceptive sensors from the agent
362 and objects. The performance of MTS3 which relies on “reconstruction loss” as the objective is yet
363 to be validated on noisy high dimensional sensor inputs like Images. Image-based experiments and
364 “non-reconstruction” based losses [5] can be taken up as future work.

365 **H Broader Impact**

366 While we do not foresee any immediate negative societal impacts of our work, we do believe that
367 machines that can replicate human intelligence at some point should be able to reason at multiple
368 levels of temporal abstractions using internal world models [5]. Having intelligent agents with
369 type 2 reasoning capabilities can have both positive and negative impacts. We believe identifying
370 and mitigating the potentially harmful effects of such autonomous systems is the responsibility of
371 sovereign governments.

372 **References**

- 373 [1] Philipp Becker, Harit Pandya, Gregor Gebhardt, Cheng Zhao, C James Taylor, and Gerhard
374 Neumann. Recurrent kalman networks: Factorized inference in high-dimensional deep feature
375 spaces. In *International Conference on Machine Learning*, pages 544–552. PMLR, 2019.
- 376 [2] Christopher M Bishop. Pattern recognition. *Machine learning*, 128(9), 2006.
- 377 [3] Peter Dayan and Geoffrey E Hinton. Feudal reinforcement learning. *Advances in neural
378 information processing systems*, 5, 1992.
- 379 [4] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for
380 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 381 [5] Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27.
382 *Open Review*, 62, 2022.
- 383 [6] Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers:
384 Exploring the stationarity in time series forecasting. In Alice H. Oh, Alekh Agarwal, Danielle
385 Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*,
386 2022.
- 387 [7] Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth
388 64 words: Long-term forecasting with transformers. In *The Eleventh International Conference
389 on Learning Representations*, 2023.
- 390 [8] Vaisakh Shaj, Philipp Becker, Dieter Buchler, Harit Pandya, Niels van Duijkeren, C James
391 Taylor, Marc Hanheide, and Gerhard Neumann. Action-conditional recurrent kalman networks
392 for forward and inverse dynamics learning. *Conference On Robot Learning*, 2020.
- 393 [9] Rohit Sonker and Ashish Dutta. Adding terrain height to improve model learning for path
394 tracking on uneven terrain by a four wheel robot. *IEEE Robotics and Automation Letters*,
395 6(1):239–246, 2020.
- 396 [10] Michael Volpp, Fabian Flürenbrock, Lukas Grossberger, Christian Daniel, and Gerhard Neu-
397 mann. Bayesian context aggregation for neural processes. In *International Conference on
398 Learning Representations*, 2020.
- 399 [11] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series
400 forecasting? *arXiv preprint arXiv:2205.13504*, 2022.
- 401 [12] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai
402 Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In
403 *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 11106–11115,
404 2021.