



Figure 1: The queries $\{q\}_{t=t_0-\tau}^{t_0}$ are the latent representations obtained from the decoder’s self-attention, same-wise keys $\{k\}_{t=t_0-\kappa}^{t_0-\tau}$ are the latent representations obtained from the encoder’s self-attention. We make a multi-horizon prediction for each time step of interest from t_0 to $t_0 + \tau$ by performing the cross-attention of the queries $\{q\}_{t=t_0-\tau}^{t_0}$ and keys $\{k\}_{t=t_0-\kappa}^{t_0-\tau}$. For example the prediction y_{t_0} is obtained by performing the cross-attention of query $q_{t_0-\tau}$ and the keys $\{k\}_{t=t_0-\kappa}^{t_0-\tau}$.

1 A Supplementary Explanation on How We Obtain Predictions

2 In Section 2.4, we describe our approach for generating forecasts for future τ time steps. This involves
3 extracting the last τ values from the integration of two models: the forecasting model (referred to
4 as Z_P) and the denoising model (referred to as Z_D). Our approach follows the encoder-decoder
5 paradigm employed in transformers. However, unlike the conventional auto-regressive decoder
6 method, where predictions are generated one element at a time, we employ a different strategy.

7 In an auto-regressive decoder, the initial prediction at time step t_0 is estimated based on the latent
8 representation at time step t_0 , denoted as the query at time step t_0 , q_{t_0} , along with the latent
9 representations from time step $t_0 - \tau$ to t_0 denoted as a sequence of keys $\{k_t\}_{t=t_0-\tau}^{t_0}$. In transformers
10 with attention mechanism, the prediction \hat{y}_{t_0} is made by estimating the similarity between the query
11 q_{t_0} and the keys $\{k_t\}_{t=t_0-\tau}^{t_0}$. Subsequent predictions are generated by recursively using the previous
12 prediction as the query for the decoder. Each prediction is then added to the key sequence, expanding
13 it for subsequent predictions.

14 This approach can be time consuming due to its sequential nature. For multi-horizon forecasting, we
15 adopt a different approach. Following the self-attention mechanism in the encoder and decoder, the
16 final predictions are generated using cross-attention. In this process, the queries are derived from the
17 latent representations of the decoder, while the keys are derived from the latent representations of the
18 encoder. However, in this approach, the decoder generates predictions at time step t_0 by utilizing
19 the query at time step $t_0 - \tau$ and the keys from time step $t_0 - \kappa$ to $t_0 - \tau$. Rather than iteratively
20 providing predictions to the decoder, our objective is to make multi-horizon forecasts for each time
21 step all at once. Consequently, to forecast at time step $t_0 + 1$, we employ the observation from time
22 step $t_0 - \tau + 1$ as the query, and so forth. By adopting this method, we eliminate the iterative nature
23 of the prediction process. Each forecast is made directly, without relying on previous predictions. In
24 the main manuscript, we express this approach as $\{\hat{y}\}_{t=t_0}^{t_0+\tau} = \{z_{P_t}\}_{t=t_0-\tau}^{t_0} + \{z_{D_t}\}_{t=t_0-\tau}^{t_0}$. Here,
25 $\{z_{P_t}\}_{t=t_0-\tau}^{t_0}$ and $\{z_{D_t}\}_{t=t_0-\tau}^{t_0}$ refer to the predictions made by the decoder for future time steps
26 from t_0 to $t_0 + \tau$ with respect to the queries from time steps $t_0 - \tau$ to t_0 . On other words, given the
27 query at time step $t_0 - \tau$, we make a multi-step ahead prediction at time step τ , and given the query
28 at time step $t_0 - \tau + 1$, we make the prediction at time step $\tau + 1$. Please refer to Figure 1 for a
29 depicted representation on how predictions are obtained.