

# Supplemental Materials for “Another Vertical View: A Hierarchical Network for Heterogeneous Trajectory Prediction via Spectrums”

## 1 TRANSFORMER DETAILS

We employ the Transformer [1] as the backbone to encode trajectory spectrums and the scene context in the two proposed sub-networks. The Transformer used in the E-V<sup>2</sup>-Net has two main parts, the Transformer Encoder and the Transformer Decoder, both of which are made up of several attention layers.

**Attention Layers.** Multi-Head Attention operations are applied in each of the attention layers. Following [1], each layer’s multi-head dot product attention with  $H$  heads is calculated as:

$$\text{Attention}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \text{softmax}\left(\frac{\mathbf{q}\mathbf{k}^T}{\sqrt{d}}\right) \mathbf{v}, \quad (1)$$

$$\text{MultiHead}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \text{fc}\left(\text{concat}(\{\text{Attention}_i(\mathbf{q}, \mathbf{k}, \mathbf{v})\}_{i=1}^H)\right). \quad (2)$$

Here,  $\text{fc}()$  denotes one fully connected layer that concatenates all heads’ outputs. Query matrix  $\mathbf{q}$ , key matrix  $\mathbf{k}$ , and value matrix  $\mathbf{v}$ , are the three layer inputs. Each attention layer also contains an MLP (denoted as  $\text{MLP}_a$ ) to extract the attention features further. It contains two fully connected layers. ReLU activations are applied in the first layer. Formally, we have output feature  $\mathbf{f}_o$  of this layer:

$$\mathbf{f}_o = \text{ATT}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \text{MLP}_a(\text{MultiHead}(\mathbf{q}, \mathbf{k}, \mathbf{v})). \quad (3)$$

**Transformer Encoder.** The transformer encoder comprises several encoder layers, and each encoder layer contains an attention layer and an encoder MLP ( $\text{MLP}_e$ ). Residual connections and normalization operations are applied to prevent the network from overfitting. Let  $\mathbf{h}^{(l+1)}$  denote the output of  $l$ -th encoder layer, and  $\mathbf{h}^{(0)}$  denote the encoder’s initial input. For  $l$ -th encoder layer, the calculation of the layer output  $\mathbf{h}^{(l+1)}$  can be written as:

$$\begin{aligned} \mathbf{a}^{(l)} &= \text{ATT}(\mathbf{h}^{(l)}, \mathbf{h}^{(l)}, \mathbf{h}^{(l)}) + \mathbf{h}^{(l)}, \\ \mathbf{a}_n^{(l)} &= \text{Normalization}(\mathbf{a}^{(l)}), \\ \mathbf{c}^{(l)} &= \text{MLP}_e(\mathbf{a}_n^{(l)}) + \mathbf{a}_n^{(l)}, \\ \mathbf{h}^{(l+1)} &= \text{Normalization}(\mathbf{c}^{(l)}). \end{aligned} \quad (4)$$

**Transformer Decoder.** Similar to the Transformer encoder, the Transformer decoder comprises several decoder layers,

and each is stacked with two different attention layers. The first attention layer in the Transformer decoder focuses on the essential parts in the Transformer encoder’s outputs  $\mathbf{h}_e$  queried by the decoder’s input  $\mathbf{X}$ . The second layer is the same self-attention layer as in the encoder. Similar to Eq. 4, we have the decoder layer’s output feature  $\mathbf{h}^{(l+1)}$ :

$$\begin{aligned} \mathbf{a}^{(l)} &= \text{ATT}(\mathbf{h}^{(l)}, \mathbf{h}^{(l)}, \mathbf{h}^{(l)}) + \mathbf{h}^{(l)}, \\ \mathbf{a}_n^{(l)} &= \text{Normalization}(\mathbf{a}^{(l)}), \\ \mathbf{a}_2^{(l)} &= \text{ATT}(\mathbf{h}_e, \mathbf{h}^{(l)}, \mathbf{h}^{(l)}) + \mathbf{h}^{(l)}, \\ \mathbf{a}_{2n}^{(l)} &= \text{Normalization}(\mathbf{a}_2^{(l)}), \\ \mathbf{c}^{(l)} &= \text{MLP}_d(\mathbf{a}_{2n}^{(l)}) + \mathbf{a}_{2n}^{(l)}, \\ \mathbf{h}^{(l+1)} &= \text{Normalization}(\mathbf{c}^{(l)}). \end{aligned} \quad (5)$$

**Positional Encoding.** Before feeding agents representations or trajectory spectrums into the Transformer, we add the positional coding to inform the relative position of each timestep or frequency portion in the sequential inputs. The position coding  $\mathbf{f}_e^t$  at step  $t$  ( $1 \leq t \leq t_h$ ) is obtained by:

$$\begin{aligned} \mathbf{f}_e^t &= (f_{e0}^t, \dots, f_{ei}^t, \dots, f_{ed-1}^t) \in \mathbb{R}^d, \\ \text{where } f_{ei}^t &= \begin{cases} \sin(t/10000^{d/i}), & i \text{ is even;} \\ \cos(t/10000^{d/(i-1)}), & i \text{ is odd.} \end{cases} \end{aligned} \quad (6)$$

Then, we have the positional coding matrix  $\mathbf{f}_e$  that describes  $t_h$  steps of sequences:

$$\mathbf{f}_e = (\mathbf{f}_e^1, \mathbf{f}_e^2, \dots, \mathbf{f}_e^{t_h})^T \in \mathbb{R}^{t_h \times d}. \quad (7)$$

The final Transformer input  $\mathbf{X}_T$  is the addition of the original sequential input  $\mathbf{X}$  and the positional coding matrix  $\mathbf{f}_e$ . Formally,

$$\mathbf{X}_T = \mathbf{X} + \mathbf{f}_e \in \mathbb{R}^{t_h \times d}. \quad (8)$$

**Layer Configurations.** We employ  $L = 4$  layers of encoder-decoder structure with  $H = 8$  attention heads in each Transformer-based sub-networks. The  $\text{MLP}_e$  and the  $\text{MLP}_d$  have the same shape. Both of them consist of two fully connected layers. The first layer has 512 output units with the ReLU activation, and the second layer has 128 but does not use any activations. The output dimensions of fully connected layers used in multi-head attention layers are set to  $d = 128$ .

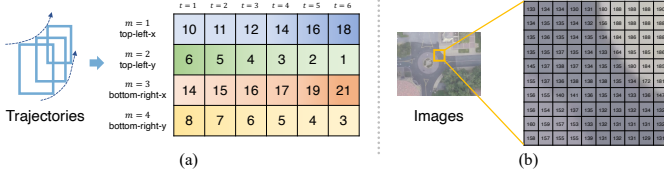


Fig. 1. Matrices views of a trajectory (2D bounding box) and an image.

## 2 LINEAR LEAST SQUARES TRAJECTORY PREDICTION

The linear least squares trajectory prediction method aims to minimize the mean square error between the predicted and agents' groundtruth trajectories. When predicting, we perform a separate least squares operation for each dimension of the  $M$ -dimensional observed trajectory  $\mathbf{X}$ . Simply, we want to find the  $\mathbf{x}_m = (b_m, w_m)^T \in \mathbb{R}^2$  ( $1 \leq m \leq M$ ), such that

$$\hat{\mathbf{Y}} = (\hat{\mathbf{Y}}_1, \hat{\mathbf{Y}}_2, \dots, \hat{\mathbf{Y}}_m, \dots, \hat{\mathbf{Y}}_M) \in \mathbb{R}^{t_f \times M},$$

$$\text{where } \hat{\mathbf{Y}}_m = \mathbf{A}_f \mathbf{x}_m = \begin{pmatrix} 1 & t_h + 1 \\ 1 & t_h + 2 \\ \dots & \dots \\ 1 & t_h + t_f \end{pmatrix} \begin{pmatrix} b_m \\ w_m \end{pmatrix}. \quad (9)$$

For one of agents' observed  $M$ -dimensional trajectory  $\mathbf{X} \in \mathbb{R}^{t_h \times M}$ , we have the trajectory slice on the  $m$ -th dimension

$$\mathbf{X}_m = (r_{m1}, r_{m2}, \dots, r_{mt_h})^T. \quad (10)$$

Suppose we have a coefficient matrix  $\mathbf{A}_h$ , where

$$\mathbf{A}_h = \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 2 & 3 & \dots & t_h \end{pmatrix}^T. \quad (11)$$

We aim to find a  $\mathbf{x}_m \in \mathbb{R}^2$ , such that the mean square  $\|\mathbf{A}_h \mathbf{x}_m - \mathbf{X}_m\|_2^2$  could reach its minimum value. Under this condition, we have

$$\mathbf{x}_m = (\mathbf{A}_h^T \mathbf{A}_h)^{-1} \mathbf{A}_h^T \mathbf{X}_m. \quad (12)$$

Then, we have the predicted  $m$ -th dimension trajectory

$$\hat{\mathbf{Y}}_m = \mathbf{A}_f \mathbf{x}_m. \quad (13)$$

The final  $M$ -dimensional predicted trajectory  $\hat{\mathbf{Y}}$  is obtained by stacking all results. Formally,

$$\hat{\mathbf{Y}} = \mathbf{A}_f (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M). \quad (14)$$

## 3 2D DFT V.S. BILINEAR STRUCTURE

We apply different transforms on *each dimension* of the trajectory to obtain the corresponding trajectory spectrums either in V<sup>2</sup>-Net or the enhanced E-V<sup>2</sup>-Net. Moreover, considering that one of our main contributions is to establish connections between trajectories (or spectrums) of different dimensions, a more natural idea might be to apply some 2D transform directly to these trajectories. However, it appears to be less effective from both theoretical analyses and experimental results. In this section, we will discuss the discrepancy between the 2D transform and the proposed bilinear structure in describing the two factors, including

the frequency response of the trajectory and the dimension-wise interactions, from different perspectives, taking DFT as an example.

**DFT on Different Directions in Trajectories.** The 2D DFT can be decomposed into two consecutive 1D DFTs performed in different directions of the target 2D matrix. The  $M$ -dimensional trajectory  $\mathbf{X} \in \mathbb{R}^{N \times M}$  is also a 2D matrix similar to 2D grayscale images. Although the 2D DFT and its variations have achieved impressive results in tasks related to image processing, they might be not directly applied to trajectories. We will analyze this problem specifically by focusing on the different directions of the transforms in the trajectory.

Fig. 1 shows an  $M = 4$  2D bounding box trajectory and an image with the matrix view. As shown in Fig. 1 (b), whether the image is sliced horizontally or vertically, the resulting vector could reflect the change in grayscale values in a particular direction. Therefore, when performing the 2D transform, the first 1D transform will extract the frequency response in a specific direction, while the second 1D transform will fuse it with the frequency response in the vertical direction.

In contrast, different slice directions of the trajectory may lead to different meanings. If the trajectories are sliced according to the time dimension, then four 1D time series will be obtained as shown in Fig. 1 (a). Applying 1D transforms to these four sequences, we can obtain four trajectory spectrums that could describe agents' frequency responses and thus describe their motions from the global plannings and interaction details at different scales. However, if the trajectory is sliced from the dimensional direction, then  $N$  (6 in the figure) 4-dimensional vectors will be obtained. These vectors contain information about agents' locations and postures at a particular moment. In addition, the focused dimension-wise interactions are also contained in these vectors. However, it should be noted that we are more interested in the relationships between the data in these vectors, *i.e.*, the "edges" between the different data. If a 1D transform is applied to these vectors, the resulting spectrum may hardly have a clear physical meaning, because the temporal or spatial adjacencies of these points are not reflected in these 4-dimensional vectors.

For example, suppose we want to apply the 1D DFT on the 4-dimensional (2D bounding box) vector  $\mathbf{x} = (x_l, y_l, x_r, y_r)^T$ . Simply, we have:

$$\begin{aligned} \mathcal{X} = \text{DFT}[\mathbf{x}] &= \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & -j & -1 & j \\ 1 & -1 & 1 & -1 \\ 1 & j & -1 & -j \end{pmatrix} \begin{pmatrix} x_l \\ y_l \\ x_r \\ y_r \end{pmatrix} \\ &= \begin{pmatrix} x_l + y_l + x_r + y_r \\ x_l - x_r - j(y_l - y_r) \\ x_l - y_l + x_r - y_r \\ x_l - x_r + j(y_l - y_r) \end{pmatrix}. \end{aligned} \quad (15)$$

Accordingly, we have its fundamental frequency portion  $\mathcal{X}[0] = x_l + y_l + x_r + y_r$  and the high-frequency portion  $\mathcal{X}[2] = x_l - y_l + x_r - y_r$ . However, since the four positions  $\{x_l, y_l, x_r, y_r\}$  do not have specific time-dependent or space-dependent like time-sequences and images, these

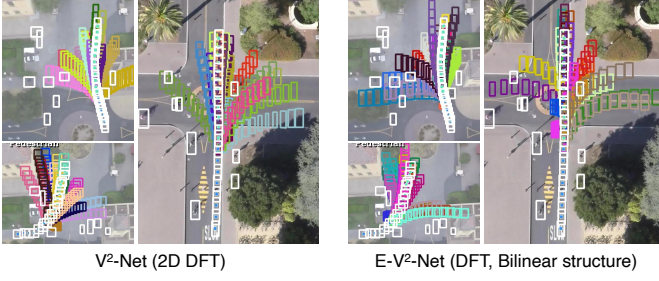


Fig. 2. Visualized comparisons of 2D DFT bilinear structure.

frequency components may hardly reflect the specific frequency response. For example, the fundamental frequencies can represent their average value for a time series, yet the values obtained by directly summing the 4 position coordinates of the 2 points of the 2D bounding box would be uninterpretable. In other words, each element in this 4-dimensional vector is relatively independent, and their connection relationships are more like a *graph* rather than a sequence where an order is required.

**Quantitative Analyses.** To verify our thoughts, we perform ablation experiments on SDD and nuScenes to compare the effects of 2D DFT and the bilinear structure quantitatively. As shown in TABLE 1, the results of APP2 and APP3 (or APP5 and APP6) show that 2D DFT does not improve quantitative trajectory prediction performance as effectively as bilinear structures. On the contrary, in the prediction of the more complex 3D bounding boxes ( $M = 6$ ), using 2D DFT instead degrades the prediction performance compared to 1D DFT when no bilinear structures are used. These experimental results validate our thoughts of not using 2D transforms but bilinear structures.

**Qualitative Analyses.** We visualize the prediction results of different models under the effect of 2D DFT and bilinear structure qualitatively. As shown in Fig. 2, the V<sup>2</sup>-Net (2D DFT) performs not as well as E-V<sup>2</sup>-Net in both the prediction of agent motions and the interactions within the bounding box. In detail, predictions given by V<sup>2</sup>-Net (2D DFT) capture fewer path possibilities in the top left prediction scenario. In addition, some predicted trajectories are with less smoothness and naturalness. For example, the prediction in color #DF6091 to the left of the bottom left prediction scene gives a turn with a large angle to observation, which could be not physical-acceptable in the actual scenario. In contrast, predictions given by E-V<sup>2</sup>-Net have not shown similar results in this scenario. On the other hand, as shown in the traffic circle prediction scenario on the right, the shape of the bounding box is not well maintained in V<sup>2</sup>-Net’s predictions, such as the prediction in color #93C3CA to turn right to across the street.

#### 4 A GRAPH-VIEW OF BILINEAR STRUCTURE

As mentioned above, the modeling of dimension-wise interactions focuses more on the relations of different trajectory dimensions, which could be difficult to represent by the 1D transform. The bilinear structure used in this manuscript

learns this connection relation through the outer-product, pooling, and fully connected networks. To make it easier to understand, we further explain it from a graph view.

Given an undirected graph  $G(t) = (V(t), E(t))$  with a time variable  $t$ , where  $V(t)$  is the set of vertices, which contains all the position information of one agent at time  $t$ , and  $E(t)$  is the set of edges, which represents the connection relationships between these vertices at time  $t$ . Formally,

$$\begin{aligned} V(t) &= \{f(r_{1t}), f(r_{2t}), \dots, f(r_{Mt})\} \\ &= \{f_{1,t}, f_{2,t}, \dots, f_{M,t}\}. \end{aligned} \quad (16)$$

where

$$f_{m,t} = f(r_{mt}) \in \mathbb{R}^d \quad (17)$$

indicates an embedding function to map these vertices into the high-dimension feature space. To establish and learn the connections between these vertices, we define the trainable adjacency matrix as:

$$A(t) = \begin{pmatrix} W_{1,1}(t) & \cdots & W_{1,M}(t) \\ \vdots & \ddots & \vdots \\ W_{M,1}(t) & \cdots & W_{M,M}(t) \end{pmatrix}. \quad (18)$$

Here, each matrix  $W_{i,j} \in \mathbb{R}^{d \times d}$  are made trainable. They are used to describe the relation between node  $i$  and node  $j$  (i.e.,  $f_{i,t}$  and  $f_{j,t}$ ).

We converge the information on the node edges by graph convolution. Formally,

$$f'_{m,t} = \sigma \left( W'_m \text{Flatten} \left( \sum_{j=1}^M W_{m,j}(t) f_{m,t} f_{j,t}^T \right) \right), \quad (19)$$

where  $\sigma$  represents a non-linear activation, and the  $W'_m$  is another trainable weight matrix. Finally, we have the refined vertices

$$V'(t) = \{f'_{1,t}, f'_{2,t}, \dots, f'_{M,t}\}. \quad (20)$$

It is worth noting that the above Eq. 19 and the bilinear structure introduced in the manuscript describe the same network structure, despite their difference in representation. To reduce unnecessary misunderstandings, we do not describe the network inference process through this graph form in the manuscript, although the use of a graph may make it easier to understand the motivation for the use of the bilinear model. In addition, all the operations above are performed on time series. If we use trajectory spectrums to replace the Eq. 17, and take the frequency variable  $n$  to instead the time variable  $t$ , we have

$$V(n) = \{f_{m,n}\}_{m=1}^M, \quad \text{where } f_{m,n} = f(s_{n,m}). \quad (21)$$

Accordingly, we have the refined vertices’ spectrum representations:

$$V'(n) = \{f'_{1,n}, f'_{2,n}, \dots, f'_{M,n}\}. \quad (22)$$

It is worth noting that in the graph view, the outer-product matrix has become the adjacency matrix of the graph  $G(n) = (V(n), E(n))$ .

#### REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in neural information processing systems*, 2017, pp. 5998–6008. 1

TABLE 1  
Validation of 2D DFT and bilinear structures with *best-of-20* on SDD (2D bounding box) the nuScenes (3D bounding box).

| No.  | Model                 | Type | T      | $N_{key}$ | BS | Dataset                       | ADE ↓ | FDE ↓ | cADE ↓ | cFDE ↓ | AIoU ↑ | FIoU ↑ |
|------|-----------------------|------|--------|-----------|----|-------------------------------|-------|-------|--------|--------|--------|--------|
| APP1 | V <sup>2</sup> -Net   | bb   | DFT    | 3         | ×  | SDD<br>(2D bounding box)      | 6.78  | 10.73 | 6.60   | 10.54  | 0.717  | 0.601  |
| APP2 | V <sup>2</sup> -Net   | bb   | 2D DFT | 3         | ×  |                               | 6.74  | 10.84 | 6.58   | 10.68  | 0.723  | 0.602  |
| APP3 | E-V <sup>2</sup> -Net | bb   | DFT    | 3         | ✓  |                               | 6.62  | 10.57 | 6.47   | 10.41  | 0.725  | 0.604  |
| APP4 | V <sup>2</sup> -Net   | 3dbb | DFT    | 2         | ×  | nuScenes<br>(3D bounding box) | 0.229 | 0.335 | 0.203  | 0.293  | 0.747  | 0.666  |
| APP5 | V <sup>2</sup> -Net   | 3dbb | 2D DFT | 2         | ×  |                               | 0.234 | 0.341 | 0.209  | 0.301  | 0.739  | 0.656  |
| APP6 | E-V <sup>2</sup> -Net | 3dbb | DFT    | 2         | ✓  |                               | 0.210 | 0.300 | 0.184  | 0.260  | 0.762  | 0.688  |