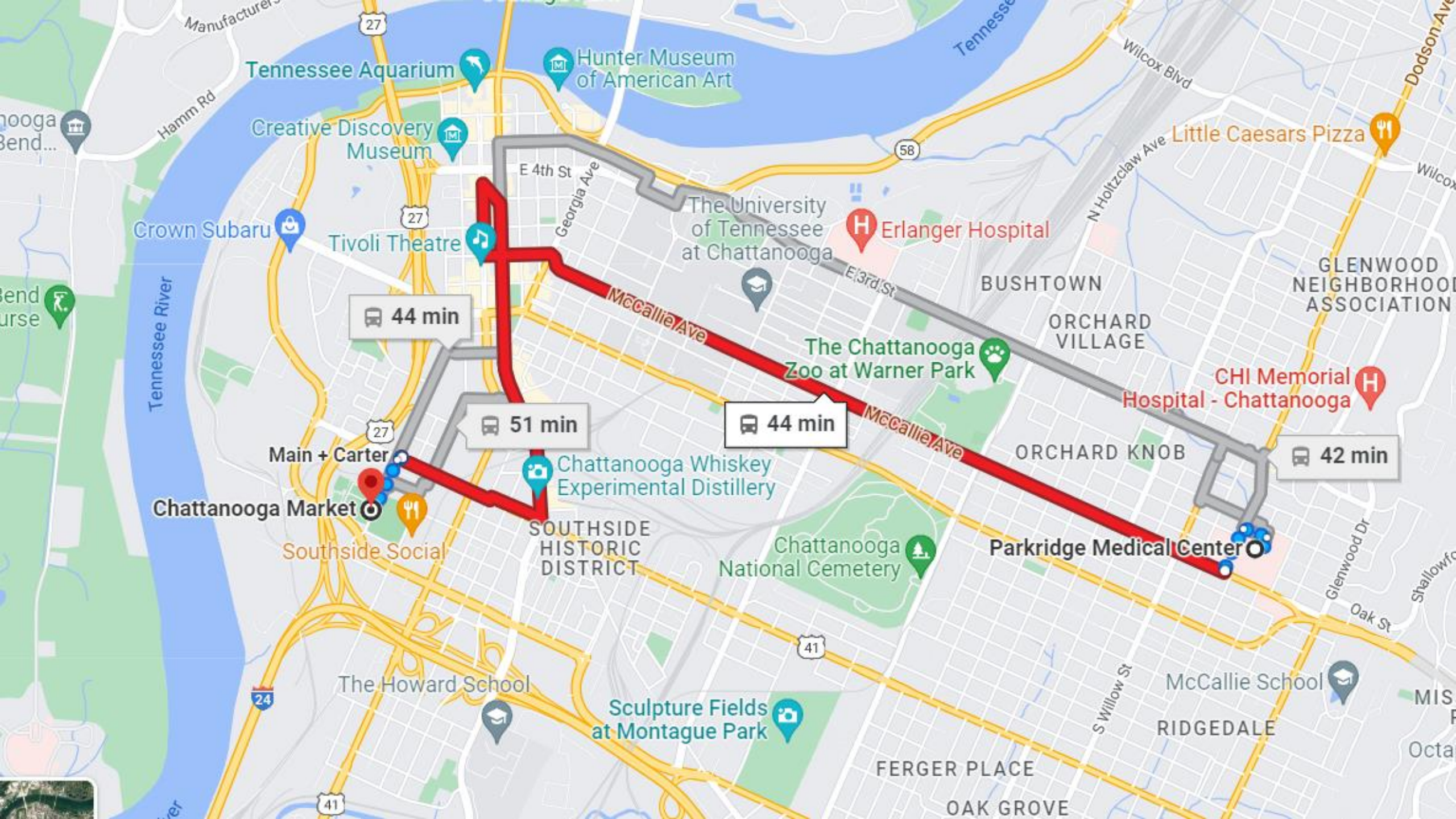


Transit Delay Time prediction using Time Series Transformers

Pravesh Koirala



CHATTANOOGA'S RIDE



Problem Statement

- Predict the delay in the Fixed-line Transit System.
- Predictors?
 - Weather (temperature, humidity, precipitation)
 - Day of the week.
 - Time of the day.
 - School breaks / holidays / Concerts or Events.
 - Nature of the route.
 - Downtown vs Suburb.

Challenges

- Nature of Predictors:

Predictor	Cardinality
Temperature	\mathbb{R}^+
Humidity	\mathbb{R}^+
Precipitation	\mathbb{R}^+
Travel Time	\mathbb{R}^+
Segment Length	\mathbb{R}^+
School Break	$\{0, 1\}$
National Holiday	$\{0, 1\}$
Day of the Week	$\{Mon, Tue, \dots Sun\}$
Time of the Day (15-min bucket)	$\{0, 1, 2, \dots 96\}$
Route Segment	$\{S_1, S_2, \dots S_{50}\}$

**There are more than 3000 route segments in the dataset, we only consider a few for the purpose of simplification*

Time Series Transformer

- An extension of the *traditional* transformer that allows you to work with continuous time predictors.
- Check <https://github.com/PraveshKoirala/Transformers-Paper> for an overview on the method.
- **EXTENSION:**
 - Change the architecture in such a way that it now works with both Categorical and Continuous predictors.

Algorithm for Time Series Forecasting

Input: $\mathbf{z}, \mathbf{x} \in \mathbb{R}^*$, two sequences of IL-Ratio inputs (+ve reals).

Output: $\mathbf{P} = \mathbb{R}^{\text{length}(\mathbf{x})}$ where $p_t \in (0,1)$ is the t -th IL-Ratio.

Hyperparameters: $l_{\max} = 10, L_{\text{enc}} = 4, L_{\text{dec}} = 4, H, d_e = d_{\text{model}}, d_{\text{mlp}} \in \mathbb{N}$

Parameters: θ includes all of the following parameters:

$\mathbf{W}_e \in \mathbb{R}^{d_e \times 1}, \mathbf{W}_p \in \mathbb{R}^{d_e \times l_{\max}}$, the token and positional embedding matrices.

For $l \in [L_{\text{enc}}]$:

| $\mathbf{W}_l^{\text{enc}}$, multi-head encoder attention parameters for layer l

| $\gamma_l^1, \beta_l^1, \gamma_l^2, \beta_l^2 \in \mathbb{R}^{d_e}$, two sets of layer-norm parameters,

| $\mathbf{W}_{\text{mlp1}}^l \in \mathbb{R}^{d_{\text{mlp}} \times d_e}, \mathbf{b}_{\text{mlp}}^l \in \mathbb{R}^{d_{\text{mlp}}}, \mathbf{W}_{\text{mlp2}}^l \in \mathbb{R}^{d_e \times d_{\text{mlp}}}, \mathbf{b}_{\text{mlp2}}^l \in \mathbb{R}^{d_e}$, MLP parameters.

For $l \in [L_{\text{dec}}]$:

| $\mathbf{W}_l^{\text{dec}}$, multi-head decoder attention parameters for layer l

| $\mathbf{W}_l^{e/d}$, multi-head cross-attention parameters for layer l .

| $\gamma_l^3, \beta_l^3, \gamma_l^4, \beta_l^4 \in \mathbb{R}^{d_e}$, two sets of layer-norm parameters,

| $\mathbf{W}_{\text{mlp3}}^l \in \mathbb{R}^{d_{\text{mlp}} \times d_e}, \mathbf{b}_{\text{mlp}}^l \in \mathbb{R}^{d_{\text{mlp}}}, \mathbf{W}_{\text{mlp4}}^l \in \mathbb{R}^{d_e \times d_{\text{mlp}}}, \mathbf{b}_{\text{mlp4}}^l \in \mathbb{R}^{d_e}$, MLP parameters.

$\mathbf{W}_u \in \mathbb{R}^{\text{length}(\mathbf{x}) \times d_e}$, the unembedding matrix.

/* Encoder portion */

```

1   $l_z \leftarrow \text{length}(\mathbf{z})$ 
2  for  $t \in [l_z]$ :  $\mathbf{e}_t \leftarrow \mathbf{W}_e[:, \mathbf{z}[t]] + \mathbf{W}_p[:, t]$ 
3   $\mathbf{Z} \leftarrow [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{l_z}]$ 
4  for  $l = 1, 2, \dots, L_{\text{enc}}$  do
5       $\mathbf{Z} \leftarrow \mathbf{Z} + \text{MHAttention}(\mathbf{Z} | \mathbf{W}_l^{\text{enc}}, \text{Mask} \equiv 1)$ 
6      for  $t \in [l_z]$ :  $\mathbf{Z}[:, t] \leftarrow \text{layer\_norm}(\mathbf{Z}[:, t] | \gamma_l^1, \beta_l^1)$ 
7       $\mathbf{Z} \leftarrow \mathbf{Z} + \mathbf{W}_{\text{mlp1}}^l \text{ReLU}(\mathbf{W}_{\text{mlp1}}^l \mathbf{Z} + \mathbf{b}_{\text{mlp1}}^l \mathbf{1}^T) + \mathbf{b}_{\text{mlp2}}^l \mathbf{1}^T$ 
8      for  $t \in [l_z]$ :  $\mathbf{Z}[:, t] \leftarrow \text{layer\_norm}(\mathbf{Z}[:, t] | \gamma_l^2, \beta_l^2)$ 
9  end
```

/* Decoder portion */

```

10  $l_x \leftarrow \text{length}(\mathbf{x})$ 
11 for  $t \in [l_x]$ :  $\mathbf{e}_t \leftarrow \mathbf{W}_e[:, \mathbf{x}[t]] + \mathbf{W}_p[:, t]$ 
12  $\mathbf{X} \leftarrow [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{l_x}]$ 
13 for  $l = 1, 2, \dots, L_{\text{dec}}$  do
14      $\mathbf{X} \leftarrow \mathbf{X} + \text{MHAttention}(\mathbf{X} | \mathbf{W}_l^{\text{dec}}, \text{Mask}[t, t'] \equiv [[t \leq t']])$ 
15     for  $t \in [l_x]$ :  $\mathbf{X}[:, t] \leftarrow \text{layer\_norm}(\mathbf{X}[:, t] | \gamma_l^3, \beta_l^3)$ 
16      $\mathbf{X} \leftarrow \mathbf{X} + \text{MHAttention}(\mathbf{X}, \mathbf{Z} | \mathbf{W}_l^{e/d}, \text{Mask} \equiv 1)$ 
17     for  $t \in [l_x]$ :  $\mathbf{X}[:, t] \leftarrow \text{layer\_norm}(\mathbf{X}[:, t] | \gamma_l^4, \beta_l^4)$ 
18      $\mathbf{X} \leftarrow \mathbf{X} + \mathbf{W}_{\text{mlp3}}^l \text{ReLU}(\mathbf{W}_{\text{mlp3}}^l \mathbf{Z} + \mathbf{b}_{\text{mlp3}}^l \mathbf{1}^T) + \mathbf{b}_{\text{mlp4}}^l \mathbf{1}^T$ 
19     for  $t \in [l_x]$ :  $\mathbf{X}[:, t] \leftarrow \text{layer\_norm}(\mathbf{X}[:, t] | \gamma_l^5, \beta_l^5)$ 
20 end
21 return  $\mathbf{P} = \mathbf{W}_u \mathbf{X}$ 
```

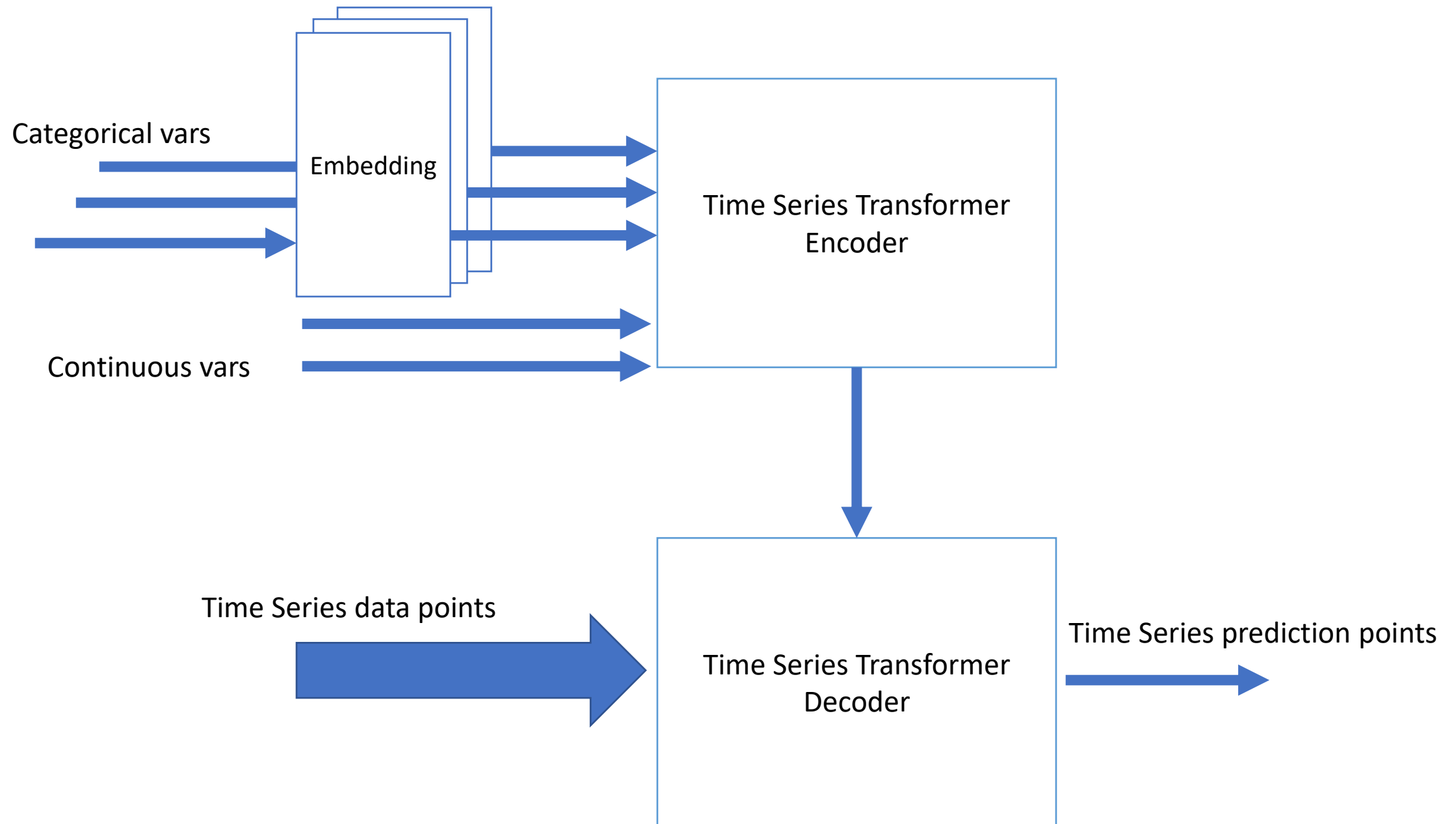
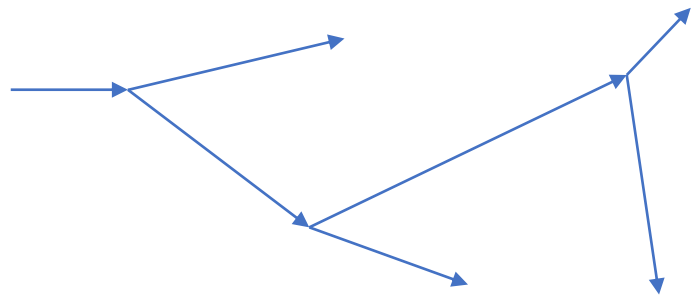


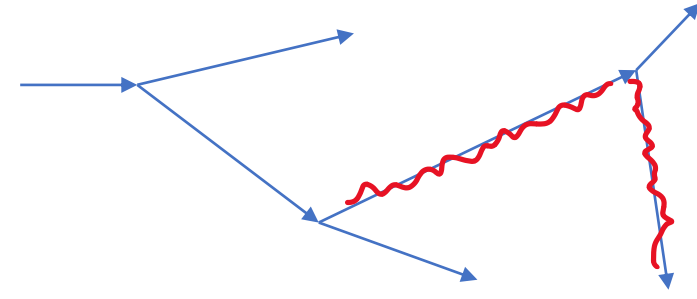
Fig: Architecture of Mixed Time Series Transformer

Critical Analysis

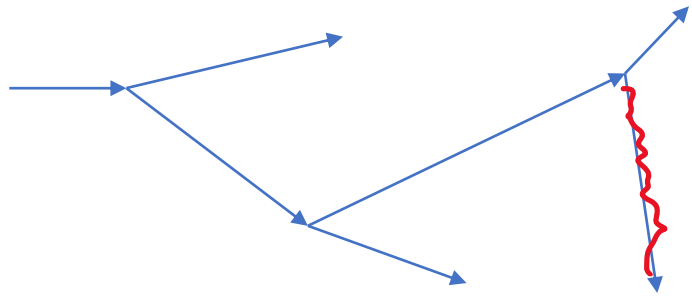
- Why not go with pure decoder architecture since it's just a next item prediction task?
 - Answer: Sticked with tried and tested model. Could potentially look into using a decoder only architecture.
- Bigger question on the inherent assumption i.e. "Delay of a segment is dependent upon past environment and past delays of the segment."
 - But realistically, delay also largely depends upon the past traffic conditions of ALL the other connected route segments. How do you encode that?



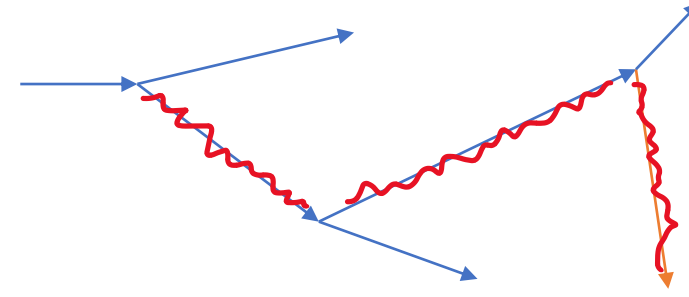
$T=0$



$T=2$



$T=1$



$T=3$

Fig: Delay propagation in Traffic Networks

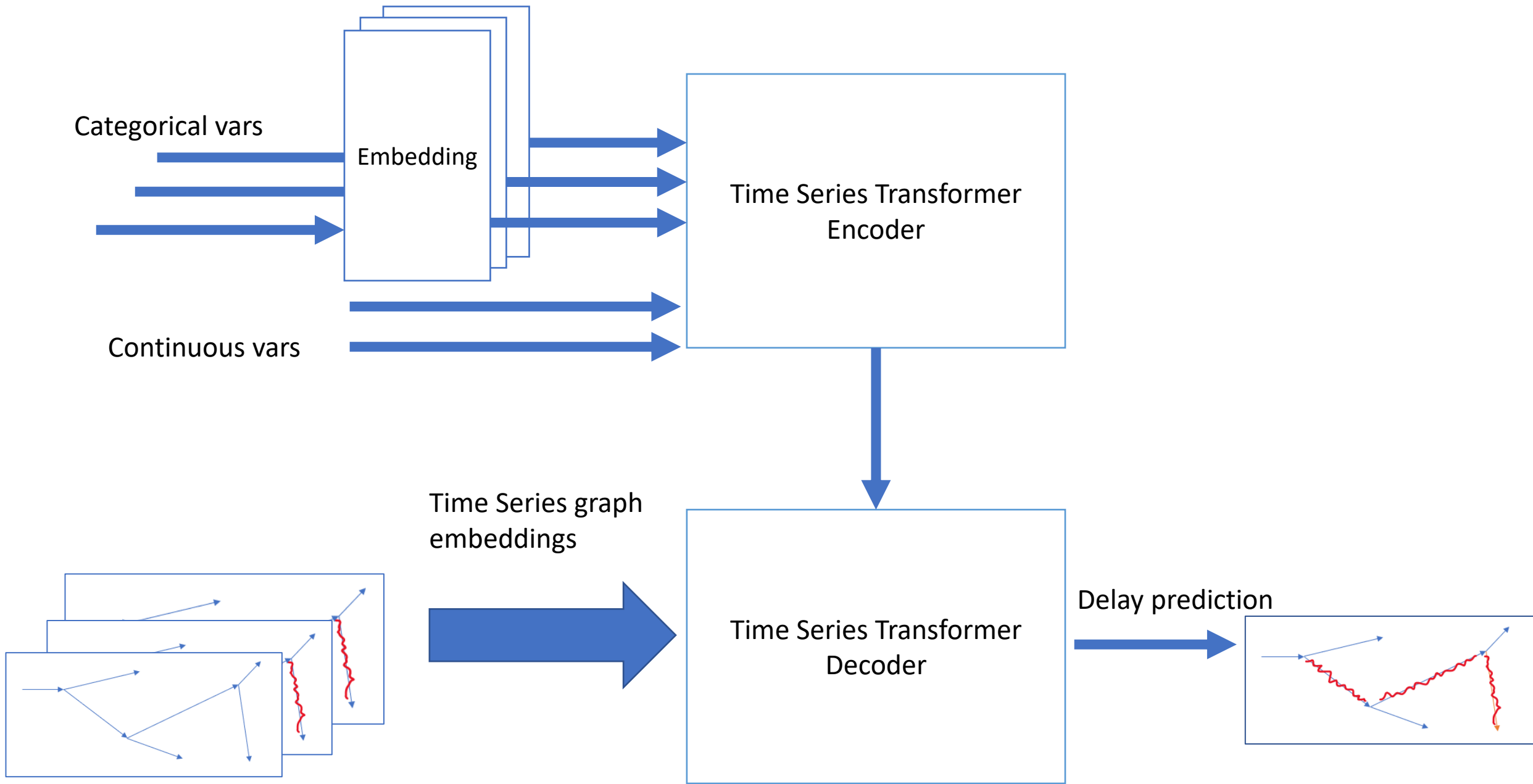


Fig: Architecture of Mixed Time Series Transformer

How do you incorporate the information of the global structure in Transformers?

- Generate Embedding using Graph Neural Networks?
- Predict \mathbb{R}^n in decoder instead of a single \mathbb{R} ?
 - Hoping that decoder self-attention will figure out a way to associate segment delays across time?