

# Sparse-Aware Spatio-Temporal Demand Forecasting with Snapshot-Based Global Attention

## 1 Input Representation

### 1.1 Static Region Features

Each region corresponds to an irregular spatial area. Static features describe functional and geographic characteristics and remain constant over time.

**Land-use composition.** Let  $C_{\text{land}}$  denote the number of land-use categories.

$$\boldsymbol{x}_r^{\text{land}} \in \mathbb{R}^{C_{\text{land}}}, \quad \sum_{k=1}^{C_{\text{land}}} x_{r,k}^{\text{land}} = 1.$$

**Example:**

$$\boldsymbol{x}_r^{\text{land}} = [0.6 \text{ (res)}, 0.3 \text{ (comm)}, 0.1 \text{ (office)}, 0, 0].$$

**POI statistics with learnable importance.** Let  $C_{\text{POI}}$  denote the number of POI categories. We use raw POI counts without explicit area normalization, since region area is separately encoded as a geographic attribute.

$$\boldsymbol{x}_r^{\text{POI}} \in \mathbb{R}_{\geq 0}^{C_{\text{POI}}}.$$

We decompose POI information into presence and magnitude components:

$$z_{r,c} \triangleq \mathbb{I}(x_{r,c}^{\text{POI}} > 0), \quad s_{r,c} \triangleq \log(1 + x_{r,c}^{\text{POI}}), \quad c = 1, \dots, C_{\text{POI}}.$$

To account for heterogeneous impact across POI categories, we introduce learnable non-negative importance weights:

$$w_c^{(z)} = \text{softplus}(\theta_c^{(z)}), \quad w_c^{(s)} = \text{softplus}(\theta_c^{(s)}).$$

The final POI feature is defined as:

$$\tilde{x}_{r,c}^{\text{POI}} = w_c^{(z)} z_{r,c} + w_c^{(s)} s_{r,c}.$$

$$\tilde{\boldsymbol{x}}_r^{\text{POI}} = [\tilde{x}_{r,1}^{\text{POI}}, \dots, \tilde{x}_{r,C_{\text{POI}}}^{\text{POI}}].$$

**Geographic attributes.** Geographic features encode region location and size:

$$\boldsymbol{x}_r^{\text{geo}} = [\text{lat}_r, \text{lon}_r, \log(1 + \text{area}_r)].$$

**Static encoder.**

$$\mathbf{u}_r^{\text{stat}} = f_{\text{stat}}(\mathbf{x}_r^{\text{land}} \|\tilde{\mathbf{x}}_r^{\text{POI}}\| \mathbf{x}_r^{\text{geo}}), \quad \mathbf{u}_r^{\text{stat}} \in \mathbb{R}^{d_s}, d_s = 32.$$

## 1.2 Temporal Context Embedding

We model calendar periodicity using learnable embedding tables.

**Day-of-week embedding.** Let  $\text{dow}(t) \in \{1, \dots, 7\}$  denote the day of week at time  $t$ . We define a learnable embedding matrix:

$$\mathbf{E}_{\text{dow}} \in \mathbb{R}^{7 \times d_t}, \quad d_t = 16.$$

The day-of-week embedding is obtained by lookup:

$$\mathbf{u}_t^{\text{dow}} = \mathbf{E}_{\text{dow}}[\text{dow}(t)].$$

**Hour-of-day embedding.** Similarly, let  $\text{hod}(t) \in \{0, \dots, 23\}$ . We define:

$$\mathbf{E}_{\text{hod}} \in \mathbb{R}^{24 \times d_t},$$

$$\mathbf{u}_t^{\text{hod}} = \mathbf{E}_{\text{hod}}[\text{hod}(t)].$$

**Holiday embedding.** We introduce a learnable holiday embedding vector:

$$\mathbf{e}^{\text{hol}} \in \mathbb{R}^{d_t}.$$

**Temporal context representation.** The final temporal context embedding is defined as:

$$\mathbf{u}_t^{\text{time}} = \mathbf{u}_t^{\text{dow}} + \mathbf{u}_t^{\text{hod}} + \mathbb{I}(\text{holiday}(t)) \cdot \mathbf{e}^{\text{hol}}.$$

## 1.3 Dynamic Demand Features (Local Lag Encoding)

We distinguish local lag encoding from the recurrent temporal horizon.

**Lag window.** Let  $\ell$  denote the local lag window length (e.g.,  $\ell = 4$ ):

$$\mathbf{y}_{r,t}^{(\ell)} = [y_{r,t-\ell+1}, \dots, y_{r,t}].$$

**Boundary handling.** For  $t - k < 0$ , we apply zero-padding and introduce missing indicators:

$$y_{r,t-k} = 0, \quad m_{r,t,k} = \mathbb{I}(t - k \geq 0).$$

**Sparsity descriptors.**

$$c_{r,t} = \sum_{i=0}^{\ell-1} \mathbb{I}(y_{r,t-i} > 0),$$

$$\Delta t_{r,t}^{\text{last}} = \min(t - \max\{\tau \leq t : y_{r,\tau} > 0\}, \Delta_{\max}).$$

The descriptor  $\Delta t_{r,t}^{\text{last}}$  captures the global recency of demand events and is conceptually independent of the local lag window length  $\ell$ . While local lag features encode short-term demand dynamics,  $\Delta t_{r,t}^{\text{last}}$  summarizes how long the region has remained inactive since its most recent non-zero demand.

Importantly,  $\Delta t_{r,t}^{\text{last}}$  is computed with respect to the *global timeline* of the dataset rather than the recurrent sequence boundary. Therefore, the last demand event may lie outside the current RNN input window. To ensure numerical stability and to reflect diminishing temporal relevance, we cap the value using a fixed threshold  $\Delta_{\max}$ . In practice,  $\Delta_{\max}$  is chosen according to dominant temporal cycles (e.g., 24 or 48 hours for daily patterns).

**Dynamic encoder.**

$$\mathbf{u}_{r,t}^{\text{dyn}} = f_{\text{dyn}}([\mathbf{y}_{r,t}^{(\ell)} \| \mathbf{m}_{r,t}^{(\ell)} \| c_{r,t} \| \Delta t_{r,t}^{\text{last}}]), \quad \mathbf{u}_{r,t}^{\text{dyn}} \in \mathbb{R}^{d_d}, d_d = 32.$$

## 2 Initial Region-Time Embedding

$$\mathbf{e}_{r,t} = \phi \left( \mathbf{W}_e [\mathbf{u}_r^{\text{stat}} \| \mathbf{u}_{r,t}^{\text{dyn}} \| \mathbf{u}_t^{\text{time}}] + \mathbf{b}_e \right), \quad \mathbf{e}_{r,t} \in \mathbb{R}^{64}.$$

## 3 Temporal State Update

We maintain a region-specific recurrent state:

$$\mathbf{h}_{r,t} = \text{GRU}(\mathbf{h}_{r,t-1}, \mathbf{e}_{r,t}), \quad \mathbf{h}_{r,0} = \mathbf{0}.$$

To combine instantaneous and history-aware representations, we apply gated fusion:

$$\mathbf{g}_{r,t} = \sigma(\mathbf{W}_g [\mathbf{e}_{r,t} \| \mathbf{h}_{r,t}] + \mathbf{b}_g),$$

$$\mathbf{s}_{r,t} = \mathbf{g}_{r,t} \odot \mathbf{h}_{r,t} + (1 - \mathbf{g}_{r,t}) \odot \mathbf{e}_{r,t}.$$

## 4 Spatial Modeling via Snapshot-Based Global Attention

At each time step, we conceptually induce a time-varying region interaction graph. Instead of explicitly constructing such graphs, we employ global self-attention as a differentiable mechanism to capture dynamic spatial dependencies.

**Attention projections.** With  $H = 4$  heads and  $d_h = 16$ :

$$\mathbf{q}_{r,t}^{(h)} = \mathbf{W}_q^{(h)} \mathbf{s}_{r,t}, \quad \mathbf{k}_{j,t}^{(h)} = \mathbf{W}_k^{(h)} \mathbf{s}_{j,t}, \quad \mathbf{v}_{j,t}^{(h)} = \mathbf{W}_v^{(h)} \mathbf{s}_{j,t}.$$

**Soft OD bias (robust normalization).** Let

$$x_{rj,t} \triangleq \log(1 + OD_{r \rightarrow j, t}).$$

We define:

$$\mu_{\tilde{OD}}(t) = \text{mean}\{x_{rj,t}\}_{r,j}, \quad \sigma_{\tilde{OD}}(t) = \text{std}\{x_{rj,t}\}_{r,j}.$$

$$\tilde{OD}_{rj,t} = \frac{x_{rj,t} - \mu_{\tilde{OD}}(t)}{\sigma_{\tilde{OD}}(t) + \epsilon}.$$

$$\lambda_{OD} = \text{softplus}(\theta_{OD}), \quad g_{OD}(r, j, t) = \lambda_{OD} \tilde{OD}_{rj,t}.$$

**Attention score.**

$$s_{rj,t}^{(h)} = \frac{(\mathbf{q}_{r,t}^{(h)})^\top \mathbf{k}_{j,t}^{(h)}}{\sqrt{d_h}} + g_{OD}(r, j, t).$$

**Aggregation.**

$$\alpha_{rj,t}^{(h)} = \frac{\exp(s_{rj,t}^{(h)})}{\sum_m \exp(s_{rm,t}^{(h)})},$$

$$\tilde{\mathbf{s}}_{r,t}^{(h)} = \sum_j \alpha_{rj,t}^{(h)} \mathbf{v}_{j,t}^{(h)}.$$

$$\tilde{\mathbf{s}}_{r,t} = \mathbf{W}_o [\tilde{\mathbf{s}}_{r,t}^{(1)} \| \cdots \| \tilde{\mathbf{s}}_{r,t}^{(H)}].$$

**Residual update.**

$$\bar{\mathbf{s}}_{r,t} = \text{LayerNorm}(\mathbf{s}_{r,t} + \tilde{\mathbf{s}}_{r,t}), \quad \epsilon = 10^{-5}.$$

## 5 Temporal Aggregation and Output

We construct training samples using a sliding temporal window of fixed length  $T$  (e.g.,  $T = 24$  hours). Specifically, for each time index  $t$ , the model consumes observations from  $\{t - T + 1, \dots, t\}$  and predicts demand at the next step  $t + 1$ . Consecutive samples are generated with a stride (up to 6 hours), resulting in overlapping sequences.

Given a temporal window  $\mathcal{T}_t = \{t - T + 1, \dots, t\}$ :

$$\mathbf{h}_{r,\tau} = \text{GRU}(\mathbf{h}_{r,\tau-1}, \bar{\mathbf{s}}_{r,\tau}).$$

**Temporal attention.**

$$\beta_{r,\tau} = \frac{\exp(\mathbf{v}_a^\top \tanh(\mathbf{W}_a \mathbf{h}_{r,\tau}))}{\sum_{k \in \mathcal{T}_t} \exp(\mathbf{v}_a^\top \tanh(\mathbf{W}_a \mathbf{h}_{r,k}))},$$

$$\bar{\mathbf{h}}_r = \sum_{\tau \in \mathcal{T}_t} \beta_{r,\tau} \mathbf{h}_{r,\tau}.$$

## 6 Sparse-Aware Output Heads

**Event probability.**

$$p_{r,t+1} = \sigma(\mathbf{w}_p^\top \bar{\mathbf{h}}_r + b_p).$$

**Conditional magnitude.**

$$\hat{y}_{r,t+1}^+ = \text{softplus}(\mathbf{w}_+^\top \bar{\mathbf{h}}_r + b_+).$$

**Final prediction.**

$$\hat{y}_{r,t+1} = p_{r,t+1} \cdot \hat{y}_{r,t+1}^+.$$

## 7 Loss Function

**Ground-truth event indicator.**

$$\delta_{r,t+1} \triangleq \mathbb{I}(y_{r,t+1} > 0).$$

**Event loss.**

$$\mathcal{L}_{\text{evt}} = - \sum_{r,t} [\delta_{r,t+1} \log p_{r,t+1} + (1 - \delta_{r,t+1}) \log(1 - p_{r,t+1})].$$

**Magnitude loss.**

$$\mathcal{L}_{\text{mag}} = \sum_{r,t} \delta_{r,t+1} \left( \frac{|\hat{y}_{r,t+1}^+ - y_{r,t+1}|}{1 + y_{r,t+1}} + \lambda |\hat{y}_{r,t+1}^+ - y_{r,t+1}| \right).$$

**OD scale regularization.**

$$\mathcal{L} \leftarrow \mathcal{L} + \eta \lambda_{OD}^2, \quad \eta = 10^{-3}.$$