

# PerDucer Formulations

2023 11033

May 2024

## 1 PerDucer Encoder Architecture

### 1.1 Behavior Tier (b-Tier) Architecture

$$\mathbf{e}_{\text{hd}}, \mathbf{e}_{\text{rel}}, \mathbf{e}_{\text{tl}} \in \mathbb{R}^d;$$

**Head-Cell**

$$\begin{aligned}\mathbf{e}'_{\text{hd}j} &= \tanh(W_{\text{hd}} \cdot \mathbf{e}_{\text{hd}j} + \mathbf{b}_{\text{hd}}); \\ \mathbf{h}'_{\text{b}j-1} &= \tanh(W_{\text{hb}} \cdot \mathbf{h}_{\text{b}j-1} + \mathbf{b}_{\text{hb}}); \\ \mathbf{h}_{\text{hd}j} &= \mathbf{e}'_{\text{hd}j} + \mathbf{h}'_{\text{b}j-1}\end{aligned}$$

**Relation-Cell**

$$\begin{aligned}\mathbf{e}'_{\text{rel}j} &= \tanh(W_{\text{rel}} \cdot \mathbf{e}_{\text{rel}j} + \mathbf{b}_{\text{rel}}); \\ \mathbf{h}'_{\text{hd}j} &= \tanh(W_{\text{hb}} \cdot \mathbf{h}_{\text{hd}j \perp (\mathbf{e}'_{\text{rel}j})} + \mathbf{b}_{\text{hb}}); \\ \mathbf{h}_{\text{rel}j} &= \mathbf{e}'_{\text{rel}j} + \mathbf{h}'_{\text{hd}j}\end{aligned}$$

**Tail-Cell**

$$\begin{aligned}\mathbf{e}'_{\text{tl}j} &= \tanh(W_{\text{tl}} \cdot \mathbf{e}_{\text{tl}j} + \mathbf{b}_{\text{tl}}); \\ \mathbf{h}'_{\text{rel}j} &= \tanh(W_{\text{hb}} \cdot \mathbf{h}_{\text{rel}j \perp (\mathbf{e}'_{\text{tl}j})} + \mathbf{b}_{\text{hb}}); \\ \mathbf{h}_{\text{b}j} &= \mathbf{e}'_{\text{tl}j} + \mathbf{h}'_{\text{rel}j}; \\ \mathbf{b}_j &= \tanh((\mathbf{p}_j + W_{\text{b}} \cdot \mathbf{h}_{\text{b}j}) + \mathbf{b}_{\text{b}})\end{aligned}$$

where: One-Hot Positional Encoding:  $\mathbf{p}_j = \mathbb{1}_j$

Where,

$$\mathbf{h}'_{\mathbf{i} \perp (\mathbf{r}_j)} = \mathbf{h}'_{\mathbf{i}} - \left( \frac{\mathbf{r}_j}{\|\mathbf{r}_j\|} \cdot \mathbf{h}'_{\mathbf{i}} \right) \cdot \left( \frac{r_j}{\|\mathbf{r}_j\|} \right)$$

## 1.2 Run Tier (r-Tier) Architecture

$$\begin{aligned}\mathbf{b}'_j &= \tanh(W_{b_r} \cdot \mathbf{b}_j + \mathbf{b}_{\mathbf{b}_r}); \\ \mathbf{h}'_{\mathbf{r}_{j-1}} &= \tanh(W_{h_r} \cdot \mathbf{h}_{\mathbf{r}_{j-1}} + \mathbf{b}_{\mathbf{h}_r}); \\ \mathbf{h}_{\mathbf{r}_j} &= \mathbf{b}'_j + \mathbf{h}'_{\mathbf{r}_{j-1}}; \\ \mathbf{r}_m &= \tanh((\mathbf{p}_m + W_r \cdot \mathbf{h}_{\mathbf{r}_j}) + \mathbf{b}_r)\end{aligned}$$

where: One-Hot Positional Encoding:  $\mathbf{p}_m = \mathbb{1}_m$

## 2 PerDucer Decoder Architecture

### 2.1 MEGA-based Encoding of Run History

Learnable damped-EMA

$$\begin{aligned} \mathbf{r}_t^{EMA} &= \boldsymbol{\alpha}_t \odot \mathbf{r}_t + (1 - \boldsymbol{\alpha}_t \odot \boldsymbol{\delta}_t) \odot \mathbf{r}_{t-1}^{EMA} \\ \text{where: } \boldsymbol{\alpha}_t &= \tanh(W_\alpha \cdot (\mathbf{r}_{t-1}^{EMA} \parallel \mathbf{r}_t) + \mathbf{b}_\alpha); \\ \boldsymbol{\delta}_t &= \tanh(W_\delta \cdot (\mathbf{r}_{t-1}^{EMA} \parallel \mathbf{r}_t) + \mathbf{b}_\delta) \end{aligned}$$

Let  $\mathbf{R}_{\tau_{\mathcal{H}}}^{EMA}$  be a matrix composed of  $\tau_{\mathcal{H}}$  rows of run  $\mathbf{r}$ . Each row is denoted by  $\mathbf{r}_t^{EMA}$ .

$$\mathbf{R}_{\tau_{\mathcal{H}}}^{EMA'} = \phi_{\text{silu}}(\mathbf{R}_{\tau_{\mathcal{H}}}^{EMA} \cdot W_{EMA} + \mathbf{b}_{EMA})$$

Contextualizing the  $\mathbf{R}_{\tau_{\mathcal{H}}}^{EMA'}$  via self-attention:

$$\begin{aligned} \mathbf{Q}_{\tau_{\mathcal{H}}}^{EMA} &= \mathbf{R}_{\tau_{\mathcal{H}}}^{EMA'} \cdot W_q + \mathbf{b}_q; \\ \mathbf{K}_{\tau_{\mathcal{H}}}^{EMA} &= \mathbf{R}_{\tau_{\mathcal{H}}}^{EMA'} \cdot W_k + \mathbf{b}_k; \\ \mathbf{V}_{\tau_{\mathcal{H}}}^{EMA} &= \mathbf{R}_{\tau_{\mathcal{H}}}^{EMA'} \cdot W_v + \mathbf{b}_v; \\ \mathbf{Z}_{\tau_{\mathcal{H}}}^{EMA} &= \text{softmax}\left(\frac{\mathbf{Q}_{\tau_{\mathcal{H}}}^{EMA} \cdot (\mathbf{K}_{\tau_{\mathcal{H}}}^{EMA})^T}{\sqrt{d_Q}}\right) \cdot \mathbf{V}_{\tau_{\mathcal{H}}}^{EMA} \end{aligned}$$

Applying forget gate to contextualized  $\mathbf{Z}_{\tau_{\mathcal{H}}}^{EMA}$ :

$$\mathbf{Z}_{\tau_{\mathcal{H}}}^{EMA^f} = \mathbf{f} \odot \mathbf{Z}_{\tau_{\mathcal{H}}}^{EMA}$$

$$\text{where: } \mathbf{f} = \phi_{\text{silu}}(\mathbf{R}_{\tau_{\mathcal{H}}}^{EMA'} \cdot W_f + \mathbf{b}_f)$$

Combining contextualized (w/ forget gate) and non-contextualized versions of  $\mathbf{R}_{\tau_{\mathcal{H}}}^{EMA'}$ :

$$\mathbf{Z}_{\tau_{\mathcal{H}}}^{EMA^C} = \phi_{\text{silu}}(\mathbf{R}_{\tau_{\mathcal{H}}}^{EMA'} \cdot W_{EMA}^C + \mathbf{Z}_{\tau_{\mathcal{H}}}^{EMA^f} \cdot W_z^C + \mathbf{b}_C)$$

Gated Residual ( $\mathbf{i}$ ) based hidden-state representation of user preference-history ( $\mathbf{R}_{\tau_{\mathcal{H}}}^h$ ):

$$\mathbf{R}_{\tau_{\mathcal{H}}}^h = \text{MEGA}(\mathbf{R}_{\tau_{\mathcal{H}}}) = \mathbf{i} \odot \mathbf{Z}_{\tau_{\mathcal{H}}}^{EMA^C} + (1 - \mathbf{i}) \odot \mathbf{R}_{\tau_{\mathcal{H}}}$$

$$\text{where } \mathbf{i} = \sigma(\mathbf{R}_{\tau_{\mathcal{H}}}^{EMA'} \cdot W_i + \mathbf{b}_i)$$

## 2.2 Prediction of upcoming runs:

$$\hat{\mathbf{r}}_{\tau_{\mathcal{F}}} = W_{\tau_{\mathcal{F}}} \cdot \mathbf{R}_{\tau_{\mathcal{H}}}^h + \mathbf{b}_{\tau_{\mathcal{F}}}; \quad \in \mathbb{R}^{|R_{KG}| \times 1}$$

$$\hat{\mathbf{p}}_{\tau_{\mathcal{F}}} = \text{sigmoid}(\hat{\mathbf{r}}_{\tau_{\mathcal{F}}}) \quad \in [0, 1]^{|R_{KG}| \times 1}$$

### 3 Task 1 Training

#### 3.1 Decoder Loss: Predicting future run sequence

$$\mathcal{L}^{\mathcal{F}}_{\mathbf{r}}(\mathbf{p}, \hat{\mathbf{p}}) = [\dots \mathcal{L}^{\mathcal{F}}_r(p_i, \hat{p}_i) \dots]$$

$$\mathcal{L}^{\mathcal{F}}_r(p_r^t, \hat{p}_r^t) = \sum_{j=1:d} \left[ p_r^{t(j)} \log \frac{1}{\hat{p}_r^{t(j)}} + (1 - p_r^{t(j)}) \log \frac{1}{(1 - \hat{p}_r^{t(j)})} \right]$$

$$\mathcal{L}^{\mathcal{F}}_r(p_r^{\tau_{\mathcal{F}}}, \hat{p}_r^{\tau_{\mathcal{F}}}) = \sum_{t=(\tau_{\mathcal{H}}+1):\tau_{\mathcal{F}}} \mathcal{L}^{\mathcal{F}}_r(p_r^t, \hat{p}_r^t)$$

#### 3.2 Run-tier (r-Tier) History Encoding Loss

$$\mathcal{L}^{\mathcal{H}}_r(p_i, \hat{p}_i) = \sum_j p_i^{(j)} \log \frac{1}{\hat{p}_i^{(j)}} + (1 - p_i^{(j)}) \log \frac{1}{(1 - \hat{p}_i^{(j)})}$$

$$\mathcal{L}^{\mathcal{H}}_{\mathbf{r}}(\mathbf{p}, \hat{\mathbf{p}}) = [\dots \mathcal{L}^{\mathcal{H}}_r(p_i, \hat{p}_i) \dots]$$

$$\begin{aligned} \nabla_{W_r}^t \mathcal{L}^{\mathcal{H}}_r(p_{\tau}, \hat{p}_{\tau}) &= \nabla_{W_r}^{t+1} \mathcal{L}^{\mathcal{H}}_r(p_{\tau}, \hat{p}_{\tau}) + \nabla_{W_r}^t \mathcal{L}^{\mathcal{H}}_r(p_t, \hat{p}_t) \\ &= \nabla_{W_r}^{t+1} \mathcal{L}^{\mathcal{H}}_r(p_{\tau}, \hat{p}_{\tau}) + [(\nabla_{h_t}^t \mathcal{L}^{\mathcal{H}}_r(p_t, \hat{p}_t) + \nabla_{W_r}^{t+1} \mathcal{L}^{\mathcal{H}}_r(p_{t+1}, \hat{p}_{t+1})) \cdot \nabla_{W_r}^t h_t] \end{aligned}$$

### 3.3 Behavior-tier (b-Tier) History Encoding Loss (BPTT)

#### 3.3.1 Local loss at time-step $t$ :

$$\mathcal{L}^{\mathcal{H}}_b(p_b^t, \hat{p}_b^t) = \sum_{j=1:d} \left[ p_b^{t(j)} \log \frac{1}{\hat{p}_b^{t(j)}} + (1 - p_b^{t(j)}) \log \frac{1}{(1 - \hat{p}_b^{t(j)})} \right]$$

$$\mathcal{L}^{\mathcal{H}}_{\mathbf{b}}(\mathbf{p}, \hat{\mathbf{p}}) = [\dots \mathcal{L}^{\mathcal{H}}_b(p_i, \hat{p}_i) \dots]$$

### 3.3.2 Loss for Back-propagation-through-time (BPTT):

**b-tier Encoder Loss:**  $\mathcal{L}^{\mathcal{H}}_b(p_b^{\tau_{\mathcal{H}}}, \hat{p}_b^{\tau_{\mathcal{H}}}) = \sum_{t=1:\tau_{\mathcal{H}}} \mathcal{L}^{\mathcal{H}}_b(p_b^t, \hat{p}_b^t)$

Decoder Loss:  $\mathcal{L}^{\mathcal{F}}_r(p_r^{\tau_{\mathcal{F}}}, \hat{p}_r^{\tau_{\mathcal{F}}})$

r-tier Cumulative Loss:  $\mathcal{L}_r(p_r^{\tau_{\mathcal{F}}+\mathcal{H}}, \hat{p}_r^{\tau_{\mathcal{F}}+\mathcal{H}})$

b-tier Cumulative Loss:  $\mathcal{L}_b(p_b^{\tau_{\mathcal{H}}+\mathcal{F}}, \hat{p}_b^{\tau_{\mathcal{H}}+\mathcal{F}}) = \mathcal{L}_r(p_r^{\tau_{\mathcal{F}}+\mathcal{H}}, \hat{p}_r^{\tau_{\mathcal{F}}+\mathcal{H}}) + \mathcal{L}^{\mathcal{H}}_b(p_b^{\tau_{\mathcal{H}}}, \hat{p}_b^{\tau_{\mathcal{H}}})$

$$\begin{aligned} \nabla_{\theta}^t \mathcal{L}_b(p_b^{\tau_{\mathcal{H}}+\mathcal{F}}, \hat{p}_b^{\tau_{\mathcal{H}}+\mathcal{F}}) &= \nabla_{\theta}^t \mathcal{L}^{\mathcal{H}}_b(p_b^{\tau_{\mathcal{H}}}, \hat{p}_b^{\tau_{\mathcal{H}}}) + \nabla_{\theta}^t \mathcal{L}_r(p_r^{\tau_{\mathcal{F}}+\mathcal{H}}, \hat{p}_r^{\tau_{\mathcal{F}}+\mathcal{H}}) \\ &= [\nabla_{\theta}^{t+1} \mathcal{L}^{\mathcal{H}}_b(p_b^{\tau_{\mathcal{H}}}, \hat{p}_b^{\tau_{\mathcal{H}}}) + \nabla_{\theta}^t \mathcal{L}^{\mathcal{H}}_b(p_b^t, \hat{p}_b^t)] + \nabla_{\theta}^t \mathcal{L}_r(p_r^{\tau_{\mathcal{F}}+\mathcal{H}}, \hat{p}_r^{\tau_{\mathcal{F}}+\mathcal{H}}) \\ &= \nabla_{\theta}^{t+1} \mathcal{L}^{\mathcal{H}}_b(p_b^{\tau_{\mathcal{H}}}, \hat{p}_b^{\tau_{\mathcal{H}}}) + \left[ \left( \nabla_{h_b^t}^t \mathcal{L}^{\mathcal{H}}_b(p_b^t, \hat{p}_b^t) + \nabla_{\theta}^{t+1} \mathcal{L}^{\mathcal{H}}_b(p_b^{t+1}, \hat{p}_b^{t+1}) \right) \cdot \nabla_{\theta}^t h_b^t \right] \\ &\quad + \nabla_{\theta}^t \mathcal{L}_r(p_r^{\tau_{\mathcal{F}}+\mathcal{H}}, \hat{p}_r^{\tau_{\mathcal{F}}+\mathcal{H}}) \end{aligned}$$

where:  $\theta \in \{W_{b_{tl}}, W_{b_{hd}}, W_{n_{hd}}, W_{n_{tl}}, W_{n_{rel}}\}$

## 4 Task 2: Next Run prediction Training

### 4.1 Query Summary Embedding Generation: Backward Translation

$$\hat{\mathbf{r}}_{\tau_{\mathcal{H}}+1}^{\mathcal{F}} = f_r(\hat{\mathbf{b}}_q, \hat{\mathbf{h}}_{r_{q-1}}) \implies \hat{\mathbf{b}}_q = f_r^{-1}(\hat{\mathbf{r}}_{\tau_{\mathcal{H}}+1}^{\mathcal{F}}, \hat{\mathbf{h}}_{r_{q-1}})$$

$$\hat{\mathbf{b}}_q = g_b(\hat{\mathbf{s}}_q, \hat{\mathbf{h}}_{summGen}) = g_b(\hat{\mathbf{s}}_q, g_{rel}(\hat{\mathbf{d}}_q, \hat{\mathbf{r}}_{summGen}))$$

$$\therefore \hat{\mathbf{h}}_{summGen} = g_{rel}(\hat{\mathbf{d}}_q, \hat{\mathbf{r}}_{summGen})$$

$$\therefore \hat{\mathbf{s}}_q = g_b^{-1}(g_{rel}(\hat{\mathbf{d}}_q, \hat{\mathbf{r}}_{summGen}), \hat{\mathbf{b}}_q) = g_b^{-1}(g_{rel}(\hat{\mathbf{d}}_q, \hat{\mathbf{r}}_{summGen}), f_r^{-1}(\hat{\mathbf{r}}_{\tau_{\mathcal{H}}+1}^{\mathcal{F}}, \hat{\mathbf{h}}_{r_{q-1}}))$$



## 5 Miscellaneous

$$\hat{\mathbf{r}}$$

$$\hat{\mathbf{r}}_{l+1}$$

$$\hat{\mathbf{r}}_{l+2}$$

$$\hat{\mathbf{r}}_{\tau}$$

$$\hat{\mathbf{p}}_{l+1}$$

$$\hat{\mathbf{p}}_{l+2}$$

$$\hat{\mathbf{p}}_{\tau}$$

$$\hat{\mathbf{P}}_{\text{next}}$$

$$\mathcal{L}^{\mathcal{H}}_b(p_b^{\tau_{\mathcal{H}}}, \hat{p}_b^{\tau_{\mathcal{H}}})$$