

# Recurrent Neural Networks for Multivariate Time Series with Missing Values

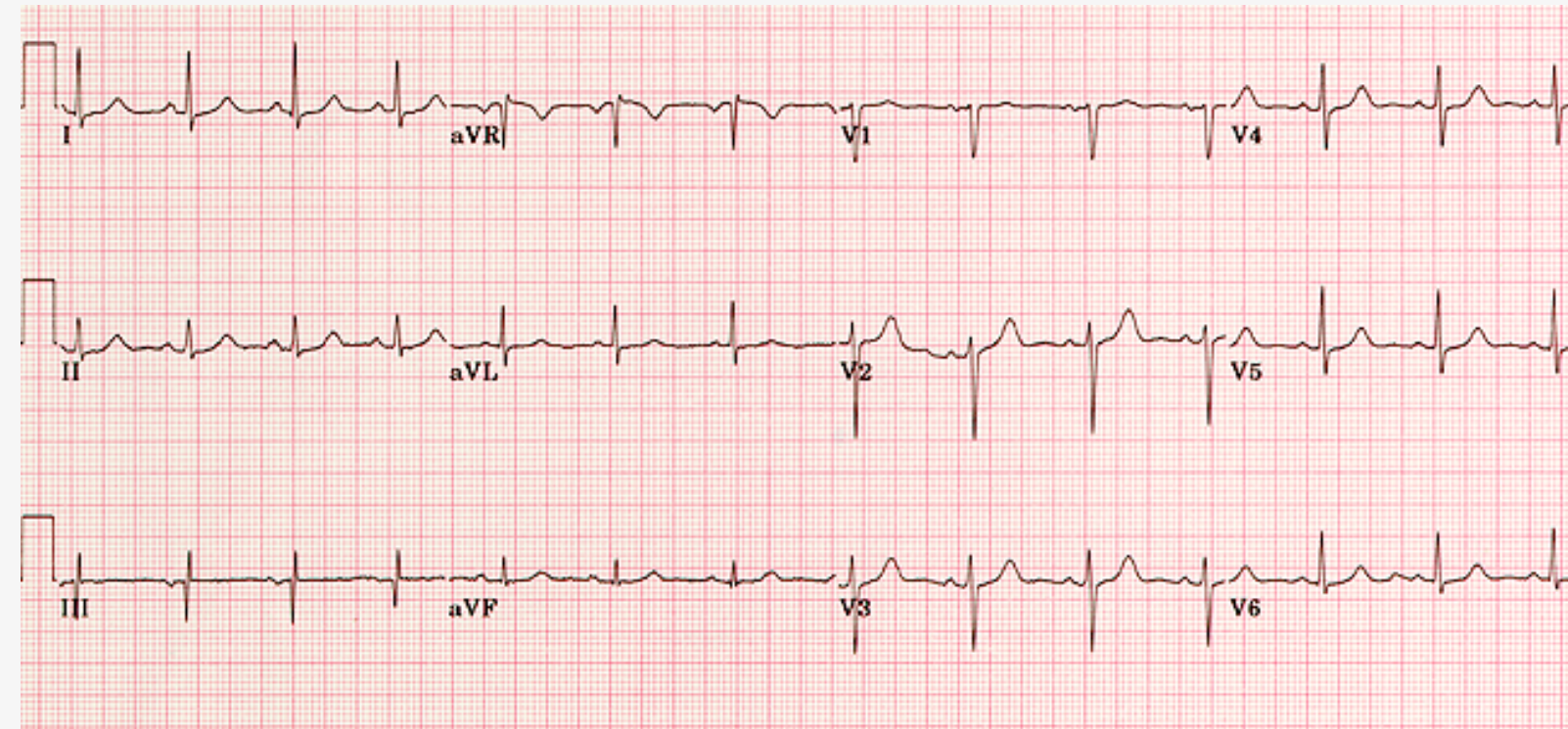
*Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag & Yan Liu*

Aslesha Pokhrel, Sujay Nagaraj

# Why do we care?

Multivariate time-series data is prevalent in a wide variety of fields:

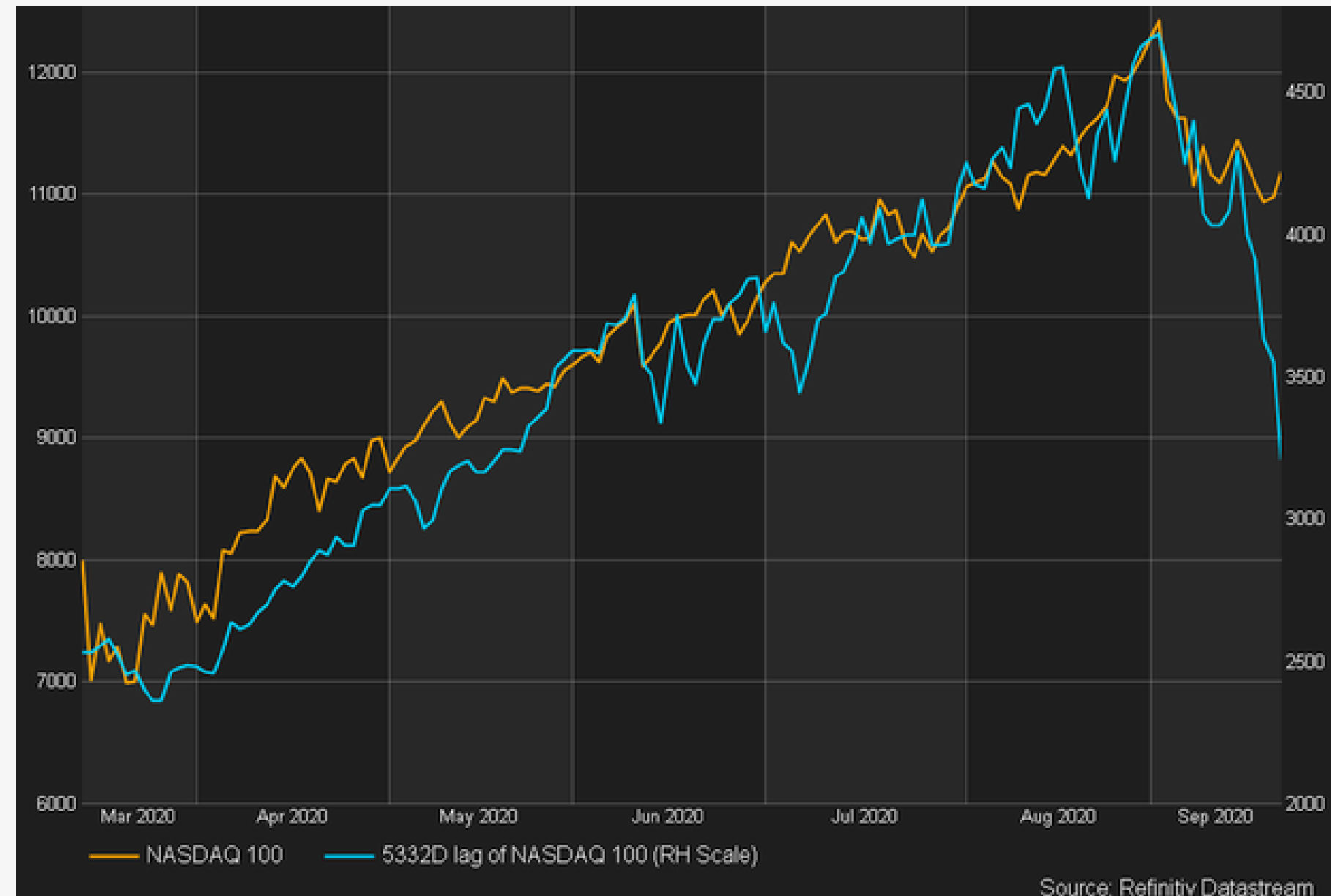
- Health care (ICU, wearables)



# Why do we care?

Multivariate time-series data is prevalent in a wide variety of fields:

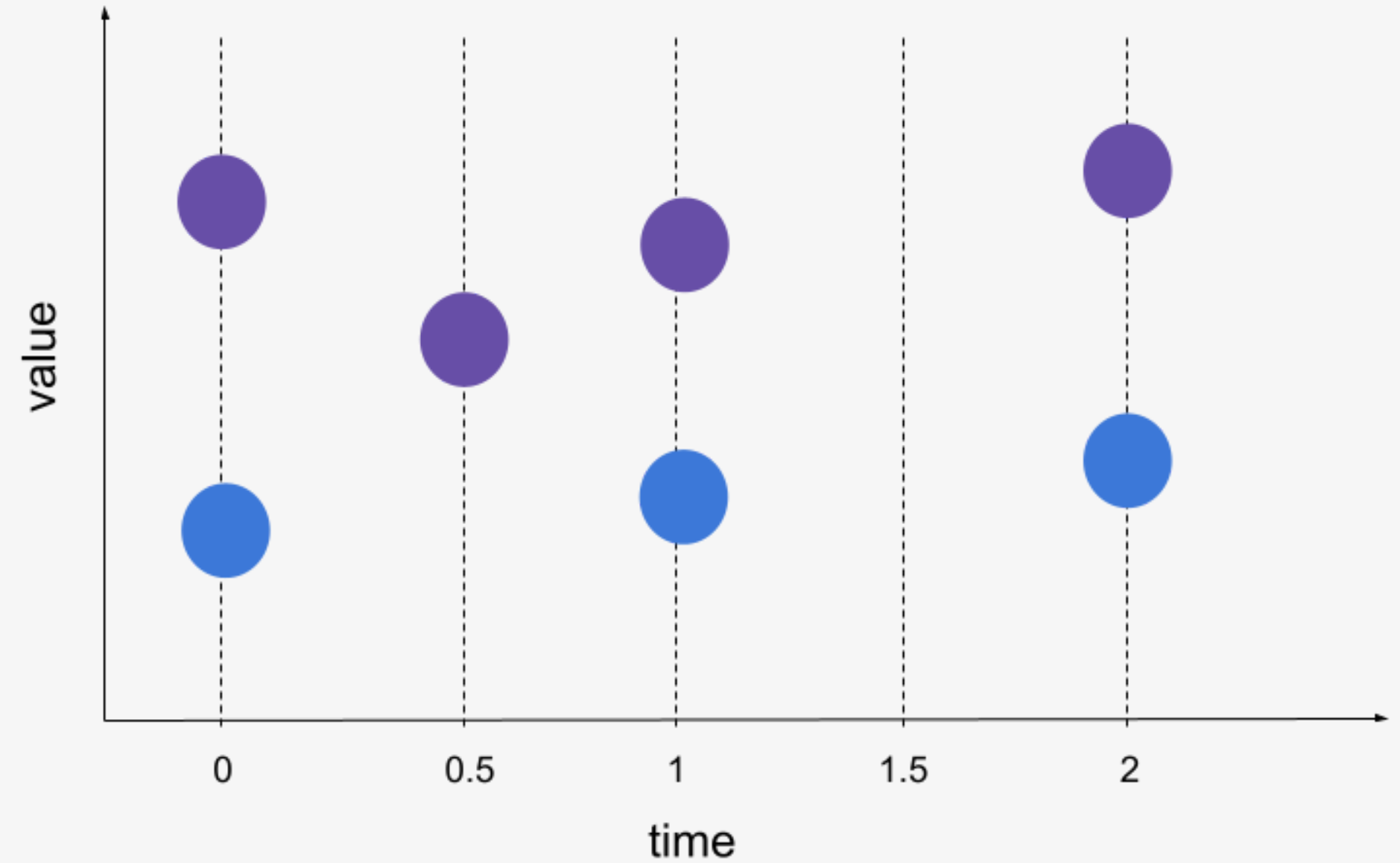
- Health care (ICU, wearables)
- Economics,
- Geoscience



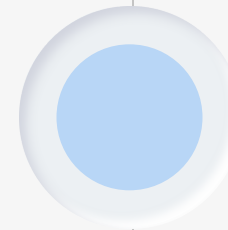
# Problem

## Missing values

- Different frequencies
- Device malfunction
- Data Corruption
- Human errors
- Intentional

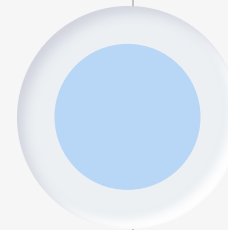


# Baseline Approaches



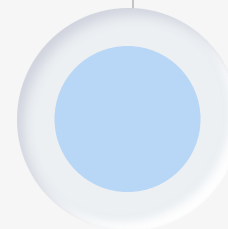
## **Omit the missing data** -

Inadequate samples when the missing rate is high.



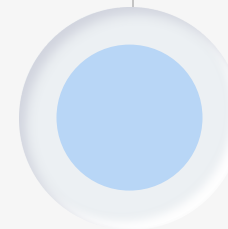
## **Data imputation** -

Methods such as smoothing, interpolation and spline are agnostic to the variable correlation and do not capture complex pattern.



## **Better methods** -

Spectral analysis, kernel methods, EM algorithm, matrix completion and matrix factorization



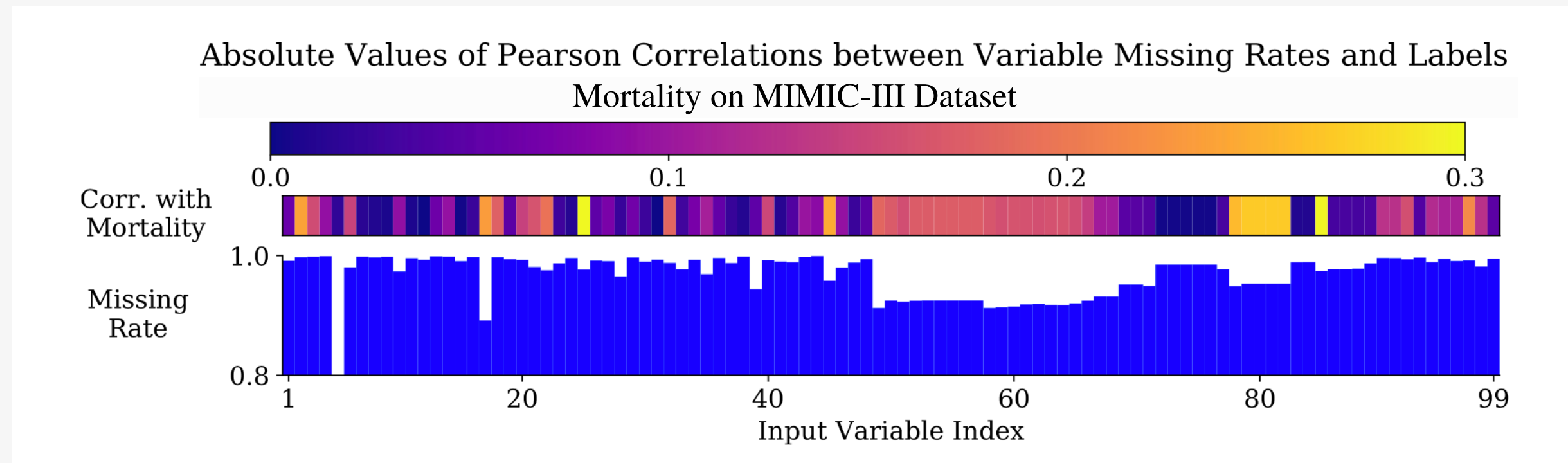
## **Multiple imputation and averaging data**

# Problem

- Assumptions of imputation method often not satisfied:
  - Small missing rate
  - Missing at random
- Results in a two-step process where imputation and prediction are separated
  - Missing patterns not effectively explored

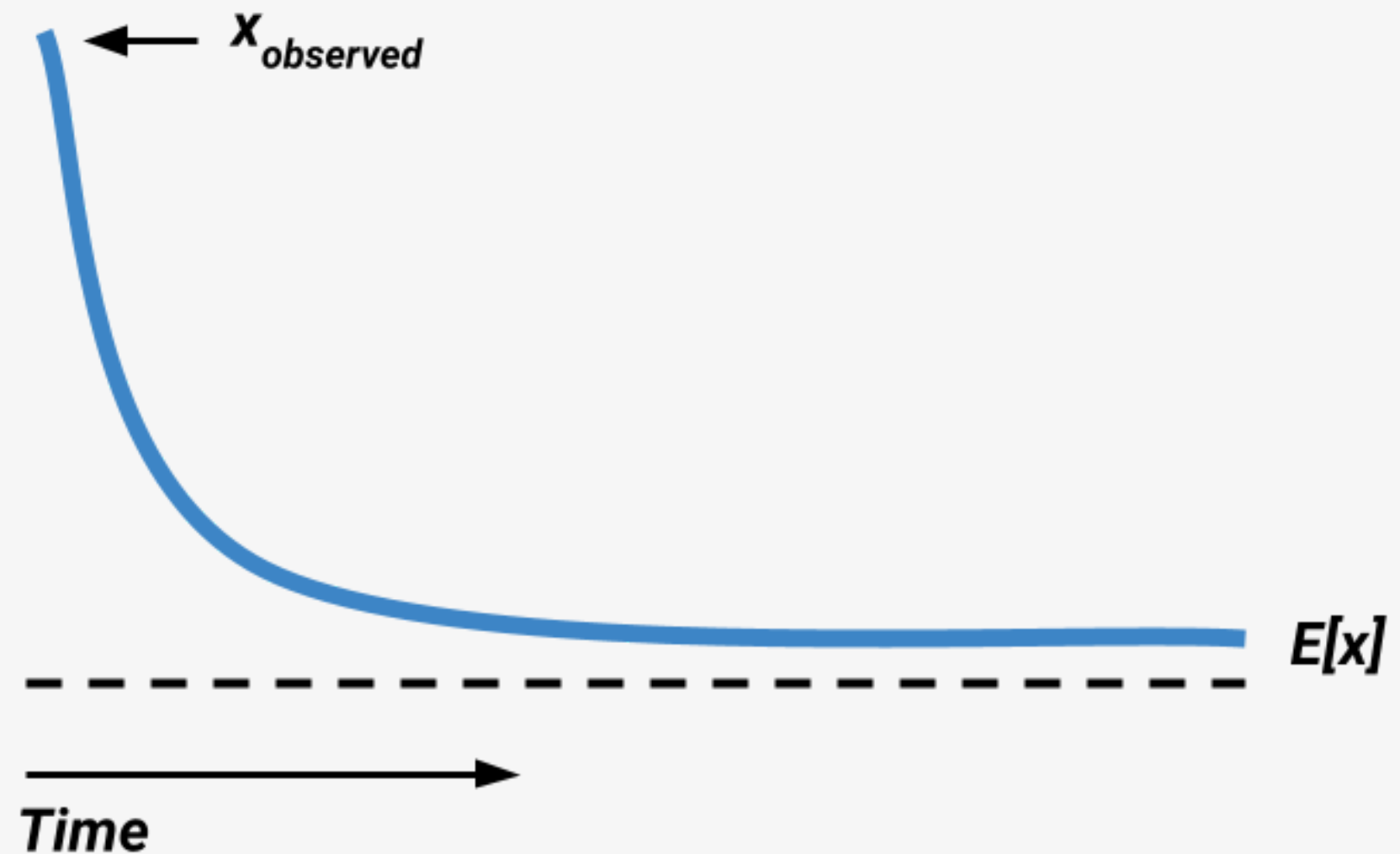
# Informative Missingness

- Missing pattern correlated with the target labels
- Provide information about labels in the supervised learning tasks



# Missingness in health data

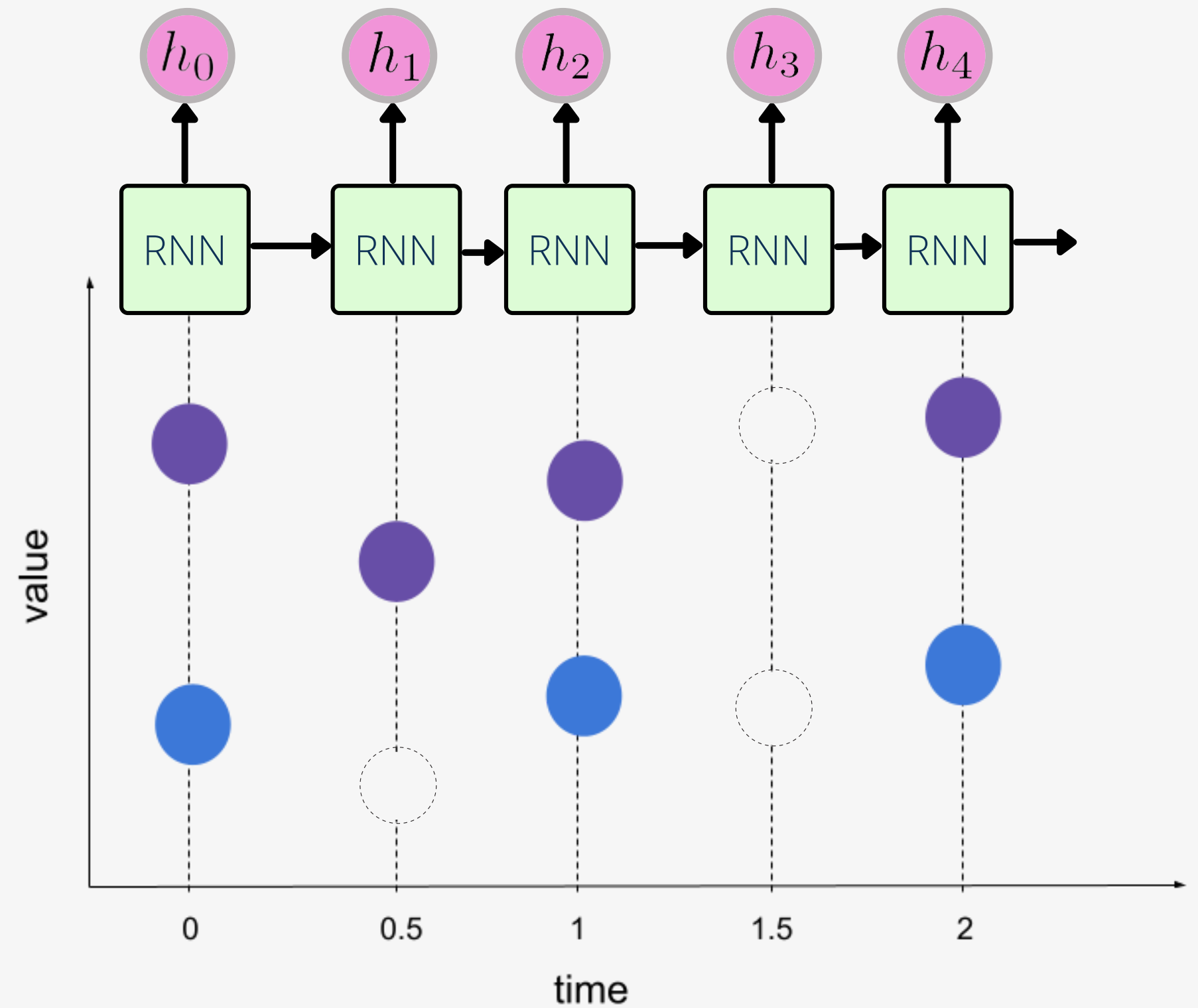
- Missing value of a variable tends to be close to some default value
  - Homeostasis
- Influence of the input variables that has been missing for a while will fade away over time



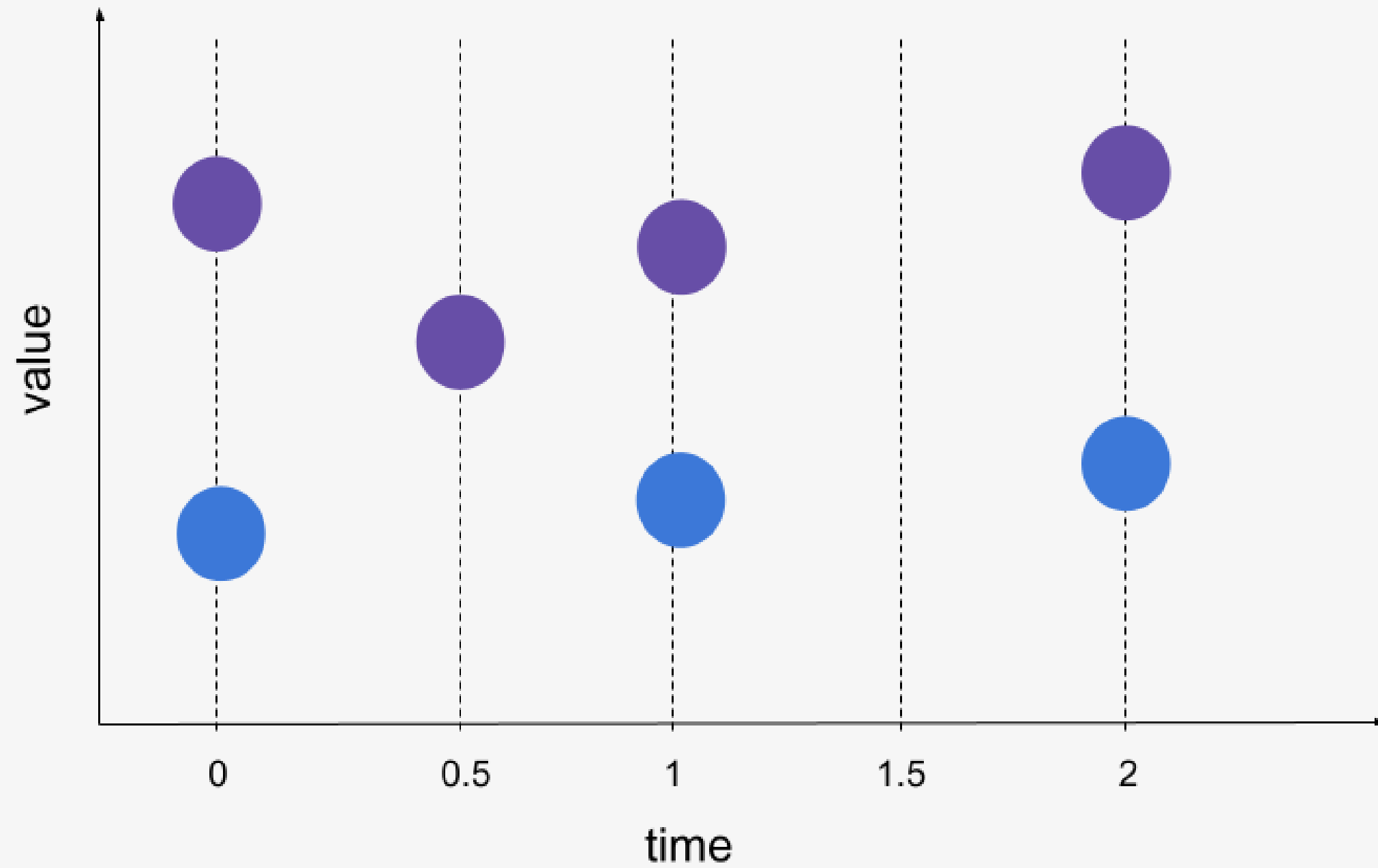


# RNN

- Models sequential data
- Auto-regressive
- Vanilla RNN suffer from vanishing gradient problem
  - GRU mitigates this problem



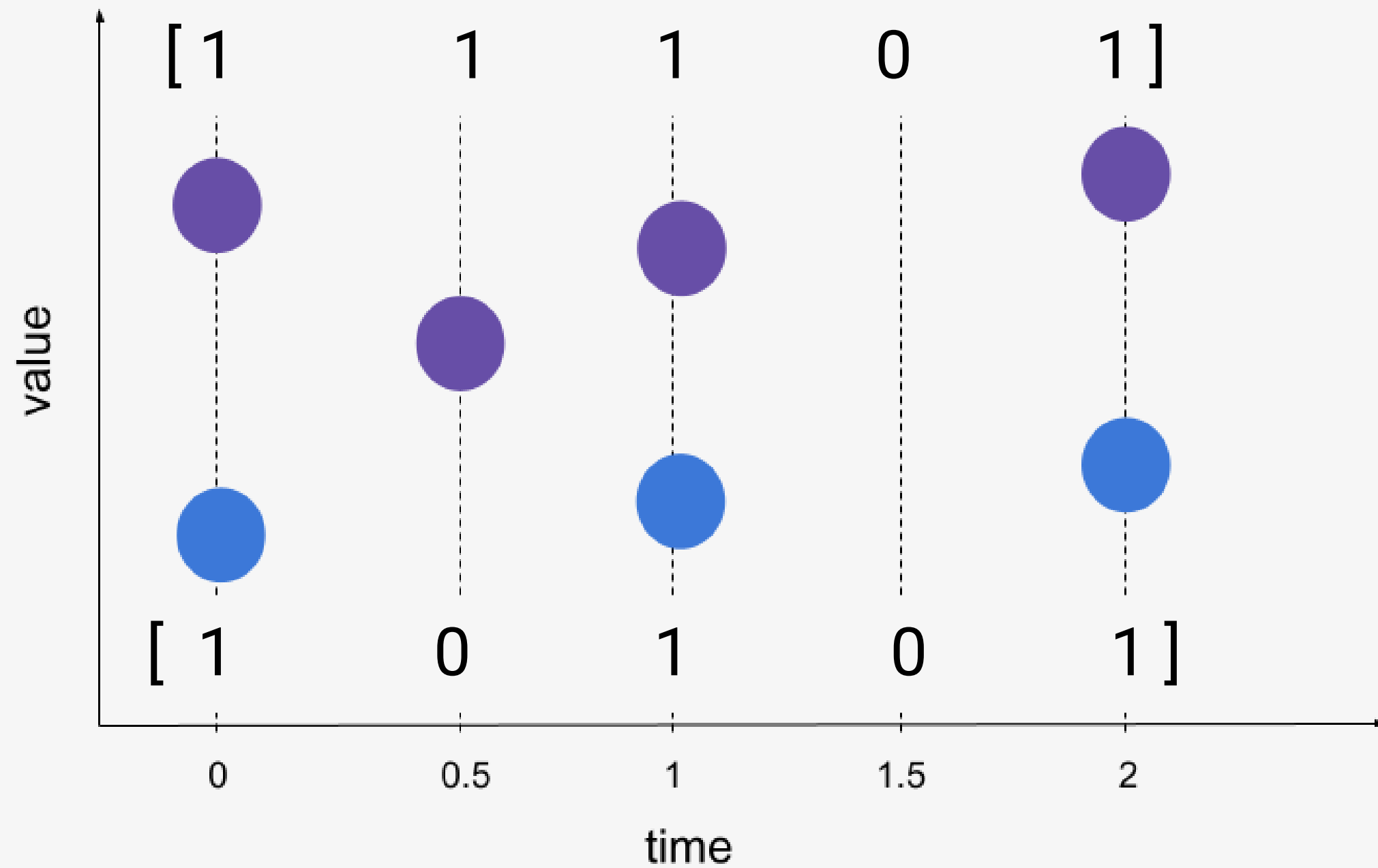
# Input



$$\mathbf{X} = \begin{bmatrix} 12 & 10 & 11 & NA & 13 \\ 6 & NA & 7 & NA & 8 \end{bmatrix}$$

$$\mathbf{s} = [0 \quad 0.5 \quad 1.0 \quad 1.5 \quad 2.0]$$

# Masking

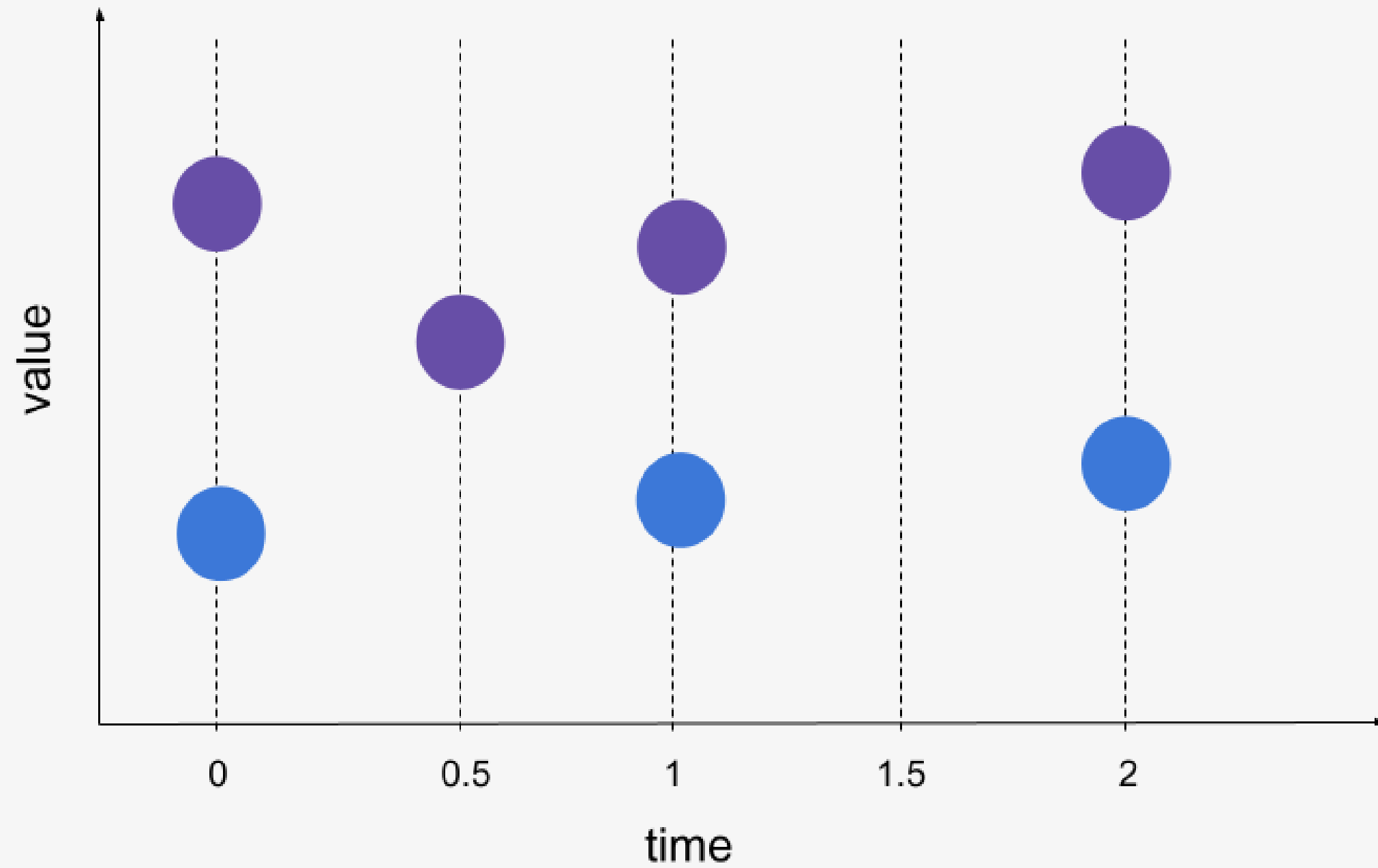


$$\mathbf{X} = \begin{bmatrix} 12 & 10 & 11 & NA & 13 \\ 6 & NA & 7 & NA & 8 \end{bmatrix}$$

$$\mathbf{s} = [0 \quad 0.5 \quad 1.0 \quad 1.5 \quad 2.0]$$

$$m_t^d = \begin{cases} 1, & \text{if } x_t^d \text{ is observed} \\ 0, & \text{otherwise} \end{cases}$$

# Masking

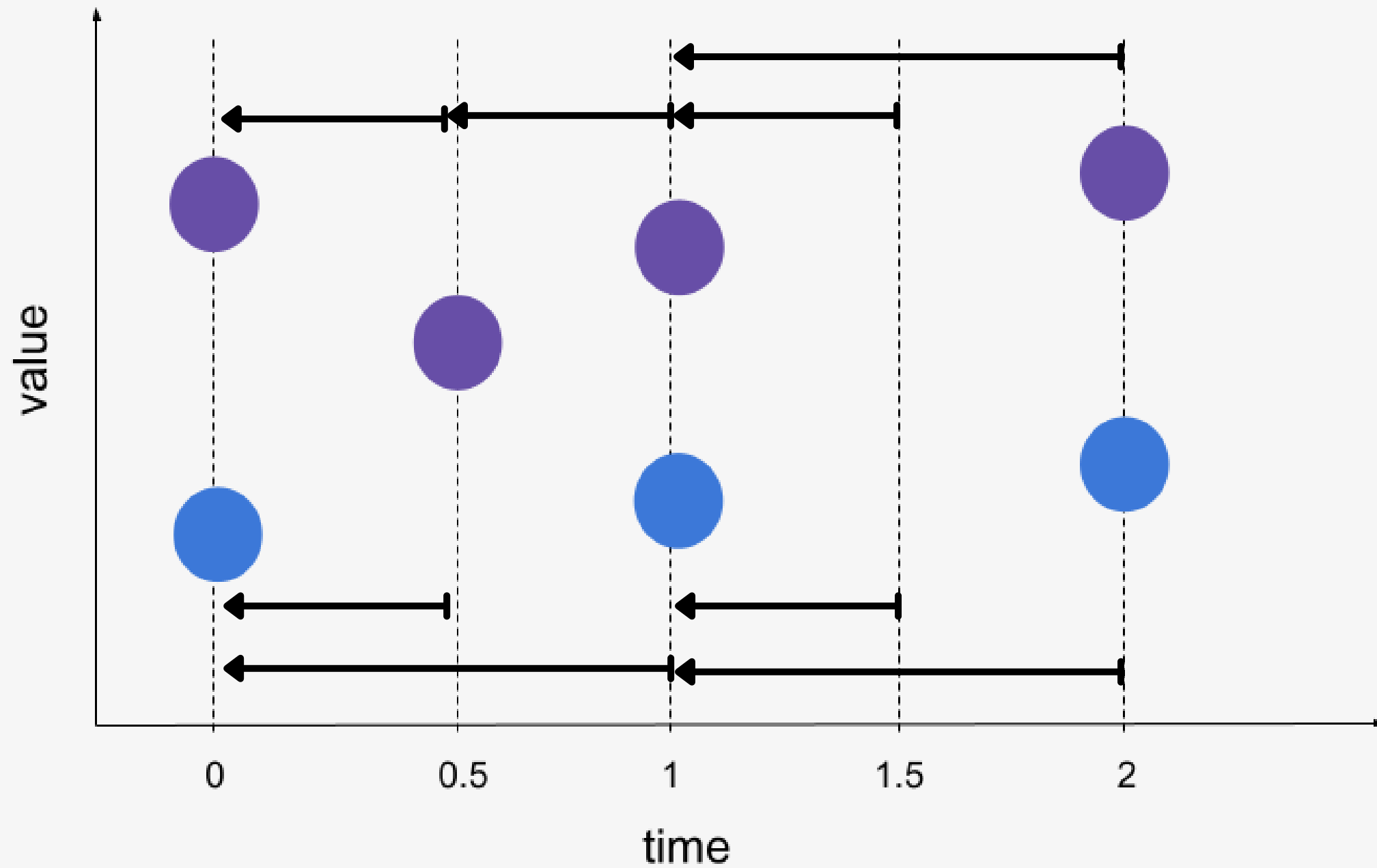


$$\mathbf{X} = \begin{bmatrix} 12 & 10 & 11 & NA & 13 \\ 6 & NA & 7 & NA & 8 \end{bmatrix}$$

$$\mathbf{s} = [0 \quad 0.5 \quad 1.0 \quad 1.5 \quad 2.0]$$

$$\mathbf{M} = \begin{bmatrix} 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$

# Time Interval



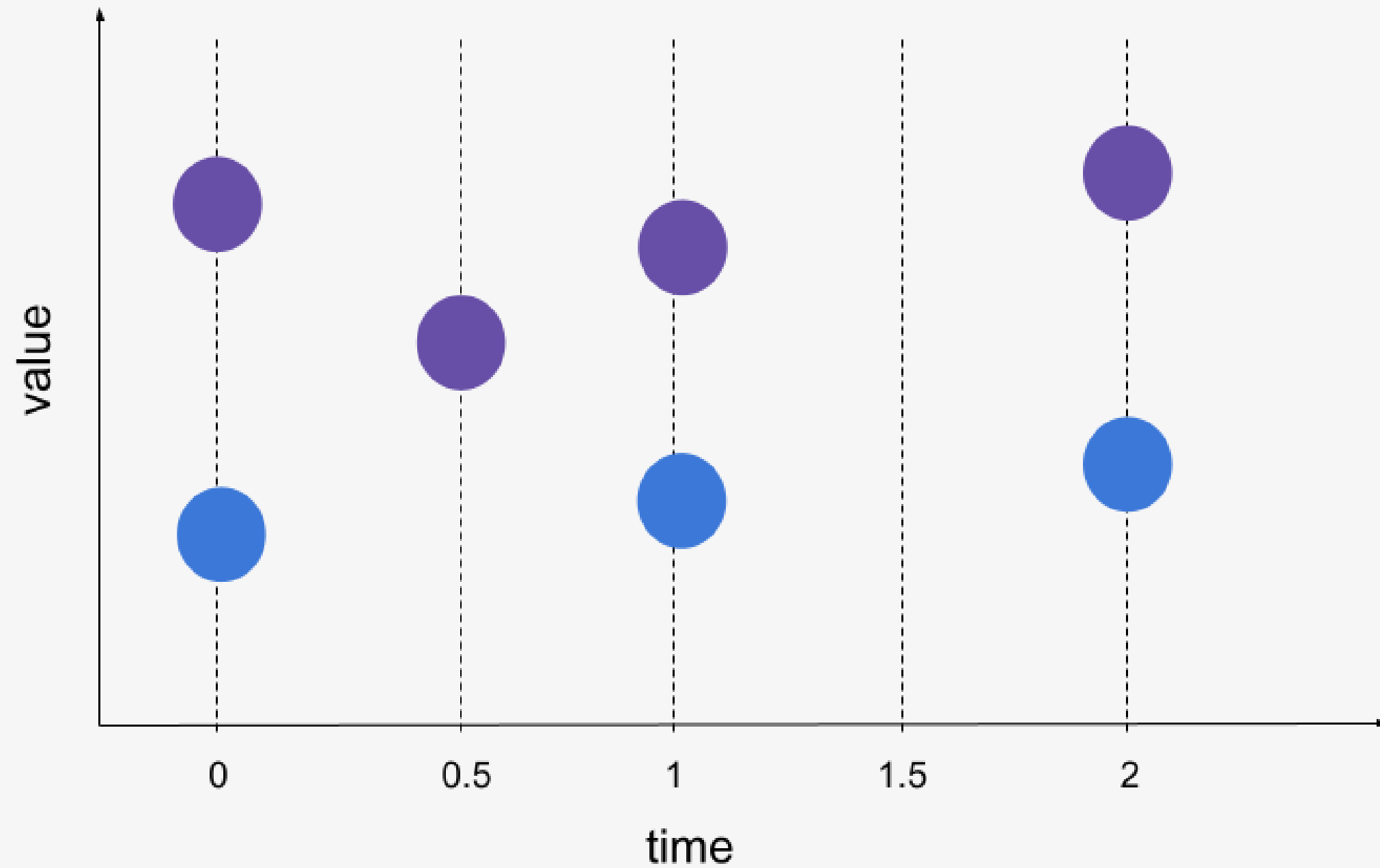
$$\mathbf{X} = \begin{bmatrix} 12 & 10 & 11 & NA & 13 \\ 6 & NA & 7 & NA & 8 \end{bmatrix}$$

$$\mathbf{s} = [0 \quad 0.5 \quad 1.0 \quad 1.5 \quad 2.0]$$

$$\mathbf{M} = \begin{bmatrix} 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$

$$\delta_t^d = \begin{cases} s_t - s_{t-1} + \delta_{t-1}^d, & t > 1, m_{t-1}^d = 0 \\ s_t - s_{t-1}, & t > 1, m_{t-1}^d = 1 \\ 0, & t = 1 \end{cases}$$

# Time Interval



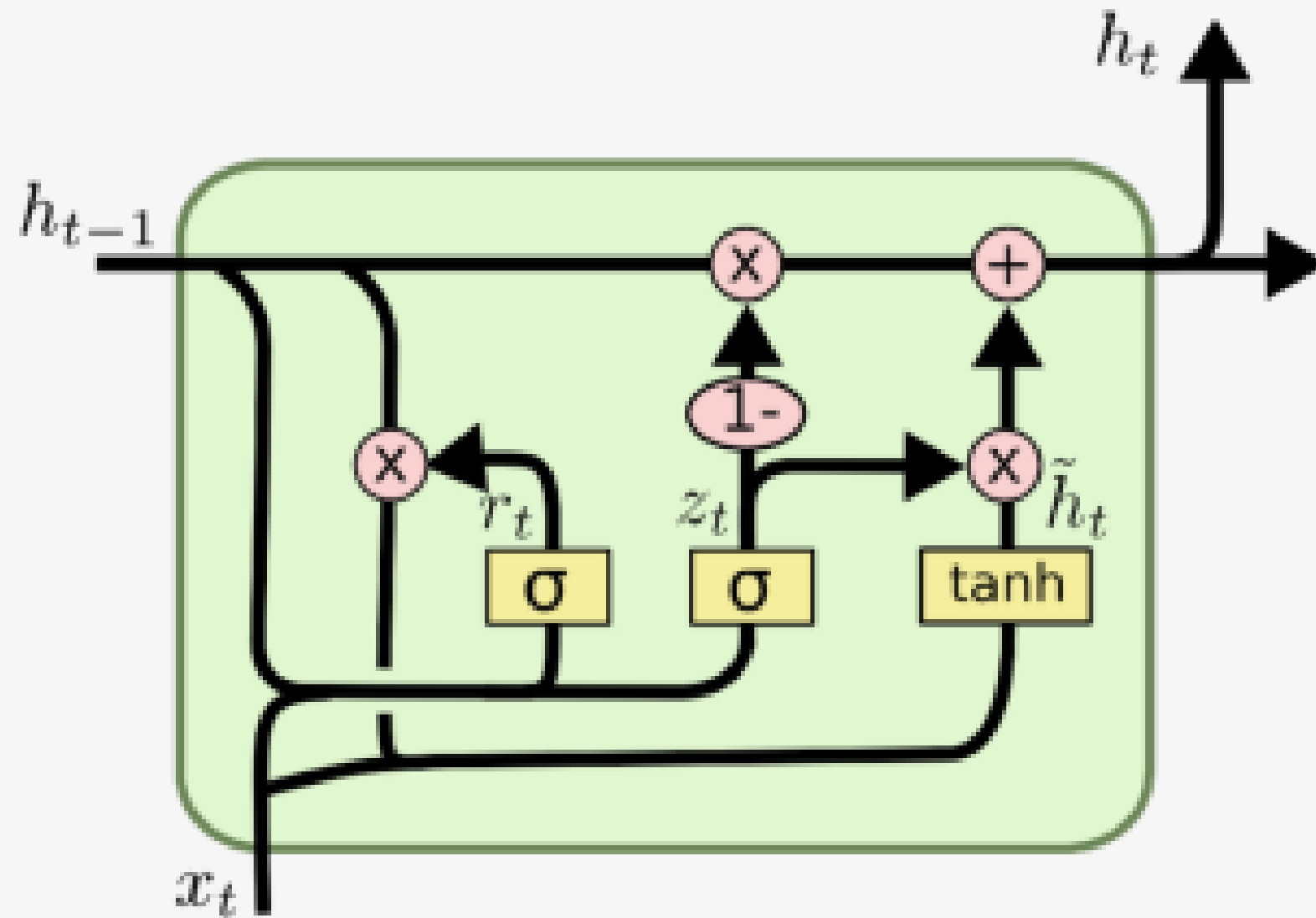
$$\mathbf{X} = \begin{bmatrix} 12 & 10 & 11 & NA & 13 \\ 6 & NA & 7 & NA & 8 \end{bmatrix}$$

$$\mathbf{s} = [0 \quad 0.5 \quad 1.0 \quad 1.5 \quad 2.0]$$

$$\mathbf{M} = \begin{bmatrix} 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$

$$\Delta = \begin{bmatrix} 0.0 & 0.5 & 0.5 & 0.5 & 1.0 \\ 0.0 & 0.5 & 1.0 & 0.5 & 1.0 \end{bmatrix}$$

# GRU



$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

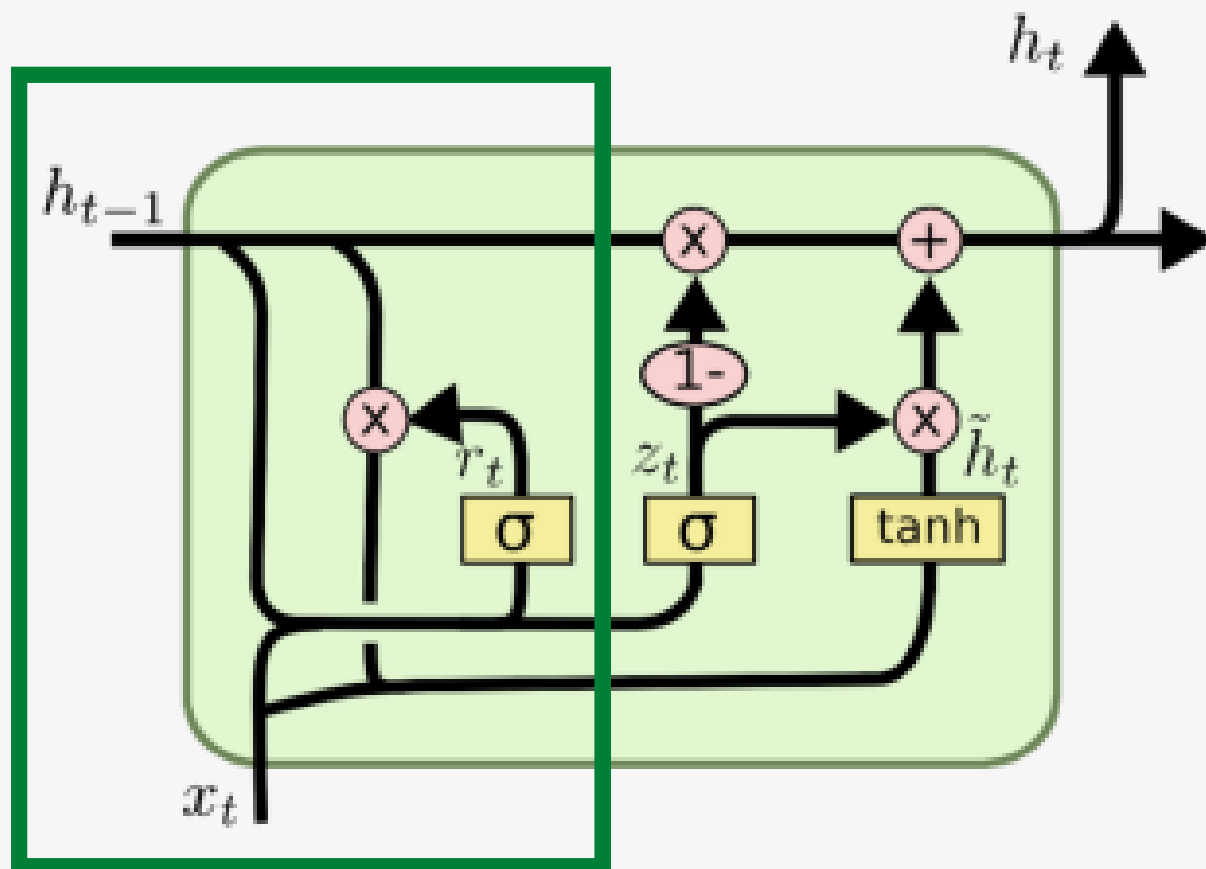
$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$$

$$\tilde{h}_t = \tanh(W x_t + U(r_t \odot h_{t-1}) + b)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

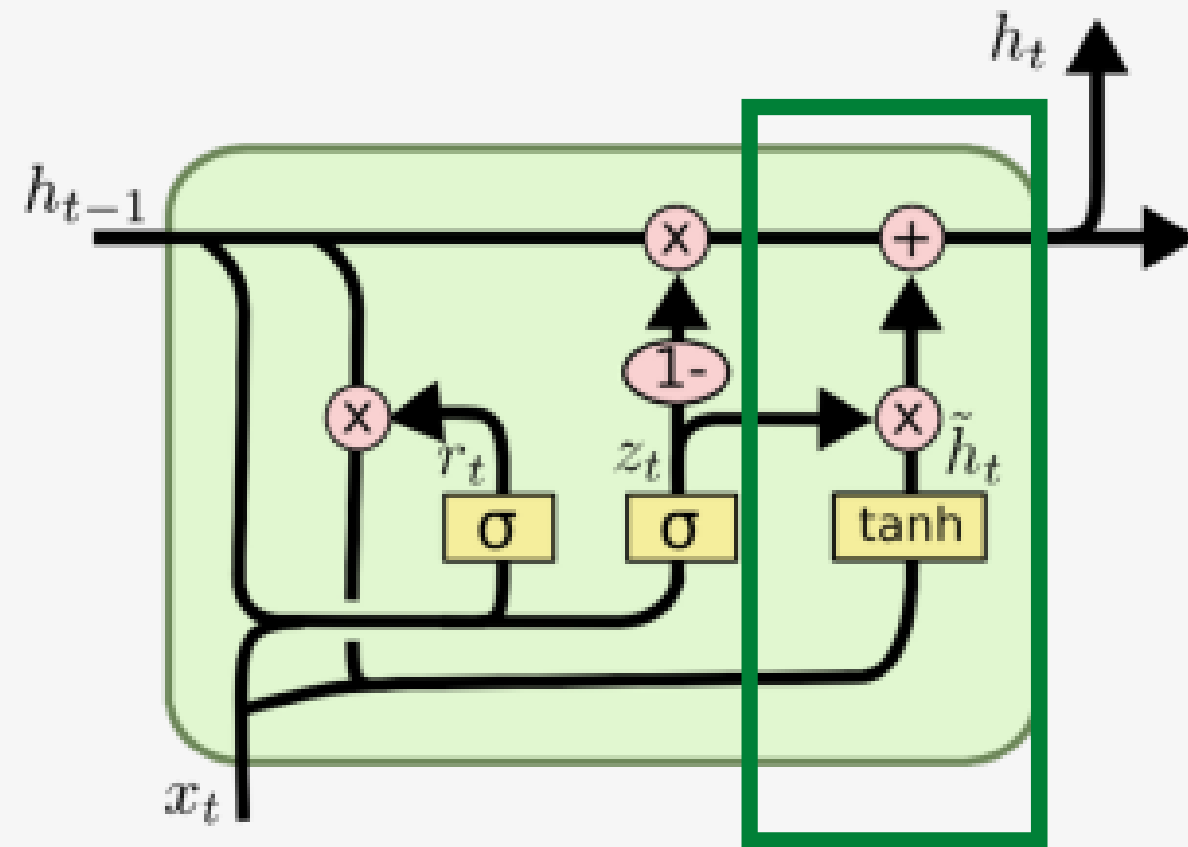
# GRU

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \quad \text{RESET GATE}$$

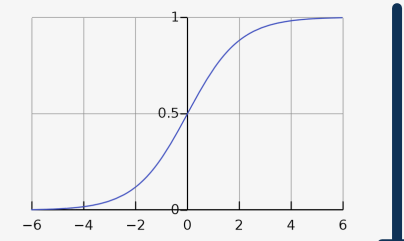




# GRU

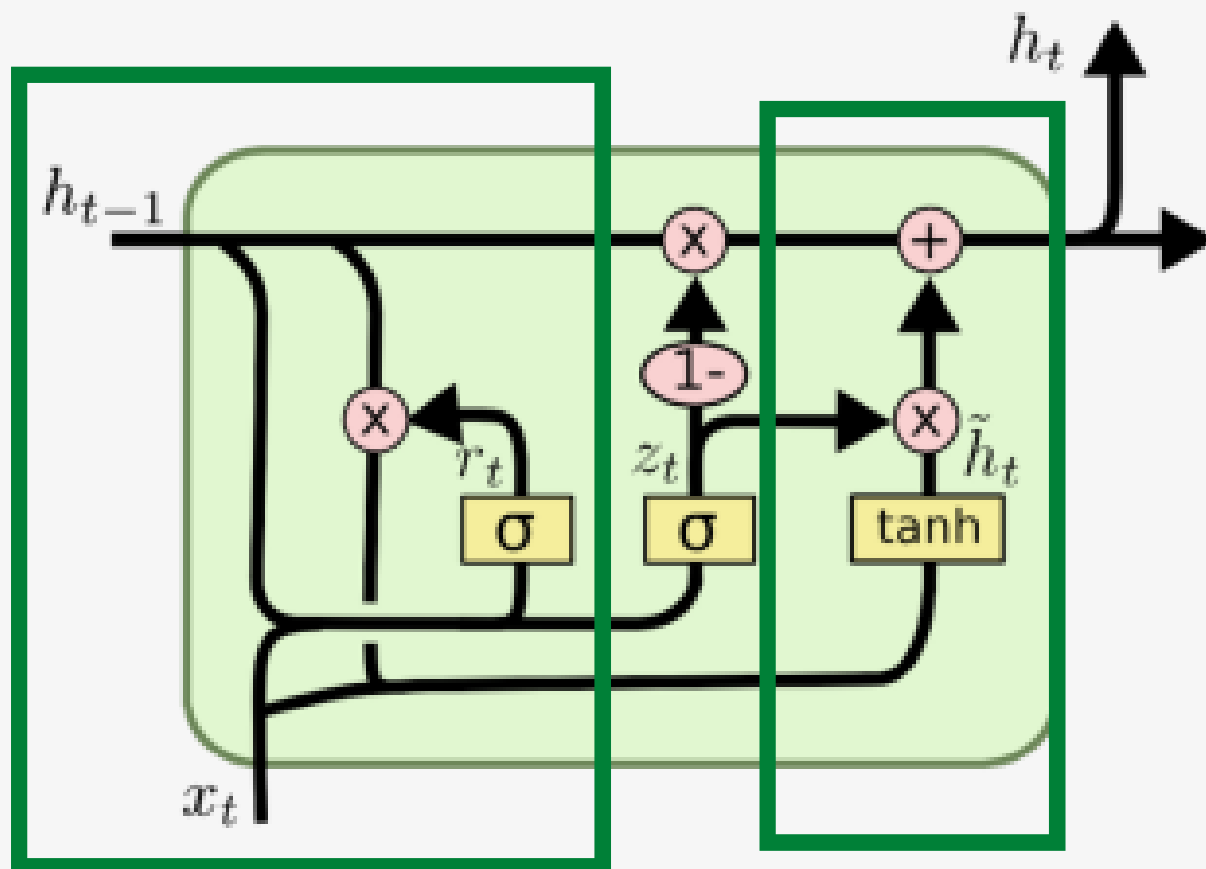


$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \quad \text{RESET GATE}$$

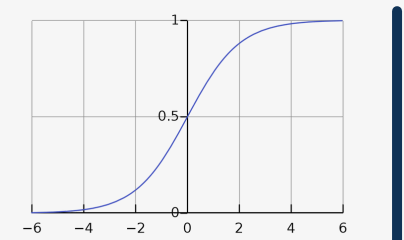


$$\tilde{h}_t = \tanh(W x_t + U(r_t \odot h_{t-1}) + b) \quad \begin{array}{l} \text{CANDIDATE} \\ \text{HIDDEN STATE} \end{array}$$

# GRU



$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \quad \text{RESET GATE}$$

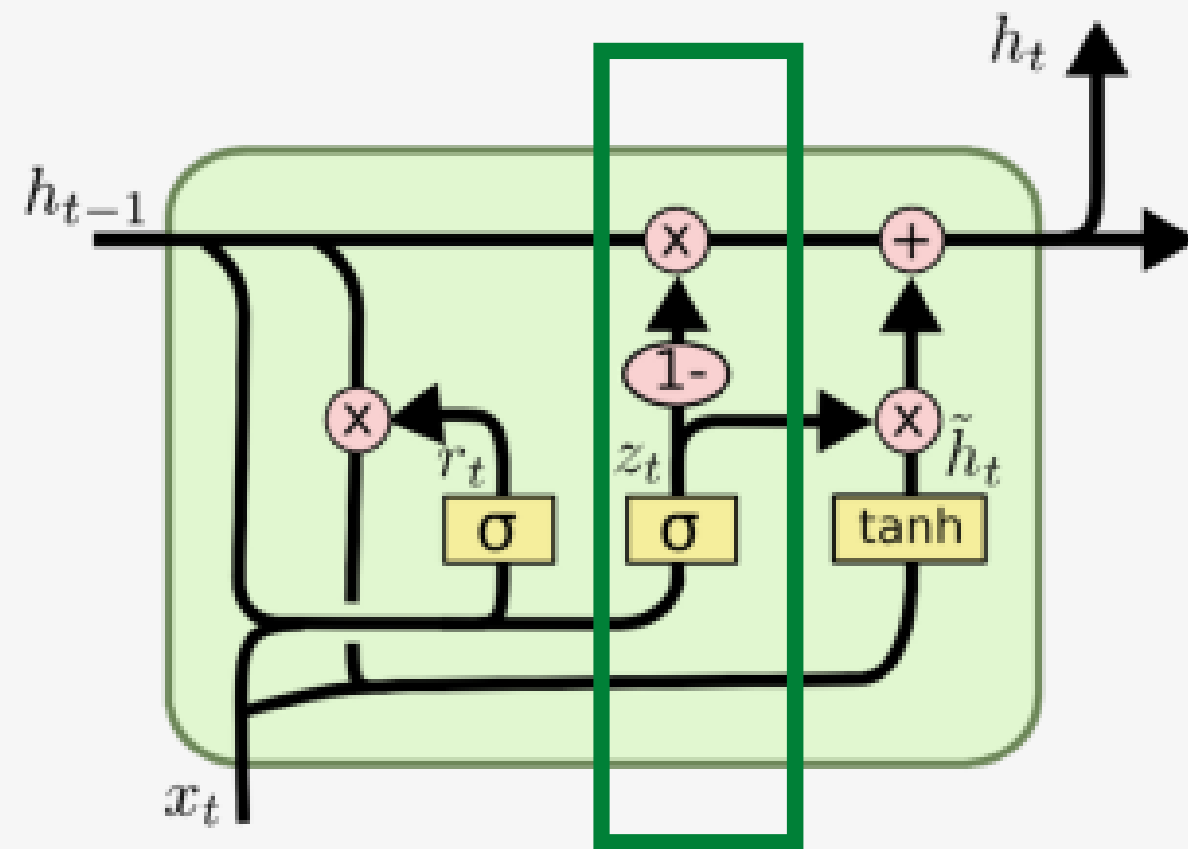


$$\tilde{h}_t = \tanh(W x_t + U(r_t \odot h_{t-1}) + b) \quad \begin{array}{l} \text{CANDIDATE} \\ \text{HIDDEN STATE} \end{array}$$

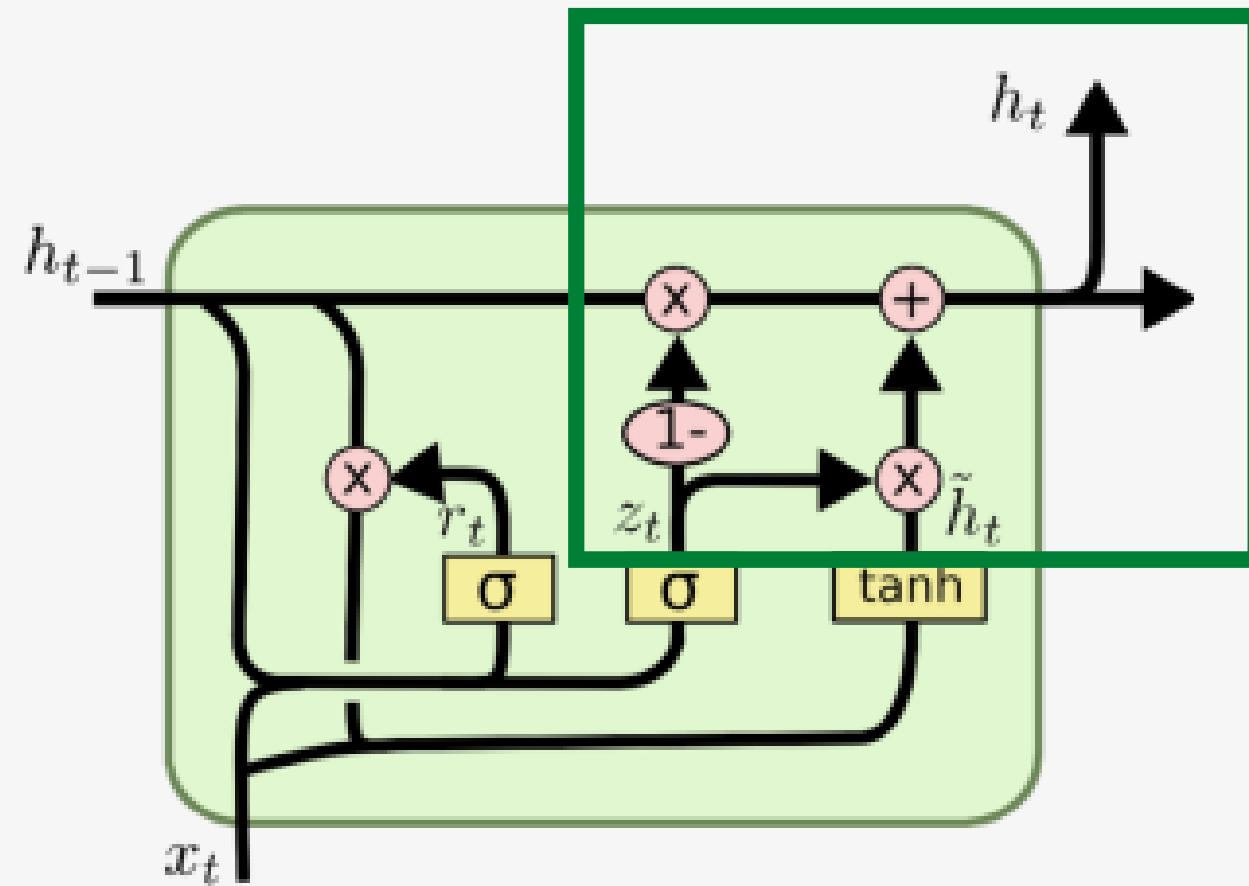
"How much do we want to remember?"

# GRU

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \text{ UPDATE GATE}$$



# GRU

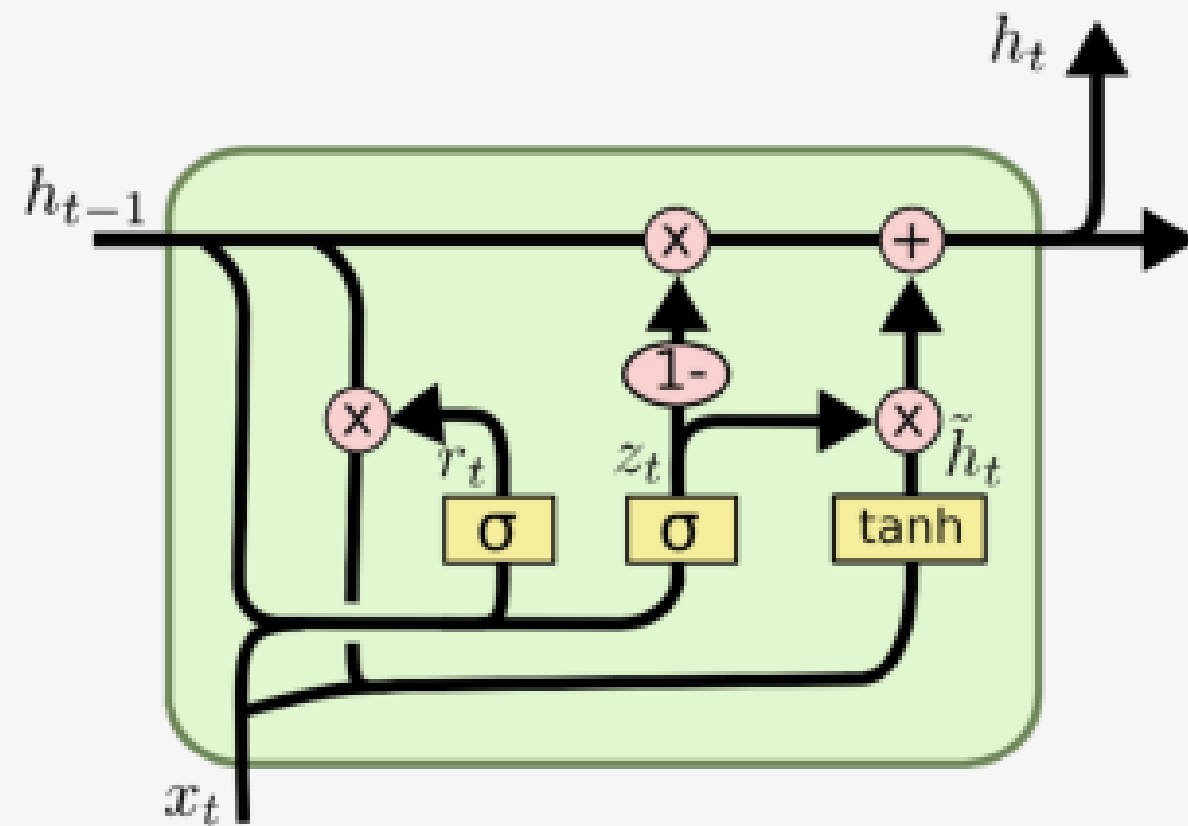


$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \text{ UPDATE GATE}$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t \text{ HIDDEN STATE}$$

CANDIDATE HIDDEN STATE

# GRU



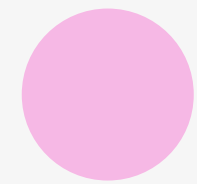
$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \text{ UPDATE GATE}$$

$$h_t = (1 - z_t) \otimes h_{t-1} + z_t \otimes \tilde{h}_t \text{ HIDDEN STATE}$$

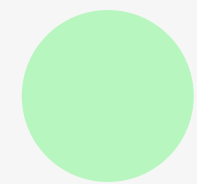
CANDIDATE HIDDEN STATE

Tradeoff between old hidden state and candidate hidden state

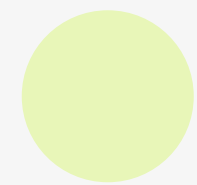
# GRU - D



**Input Decay**



**Hidden State Decay**



**Mask Vector Inputs**

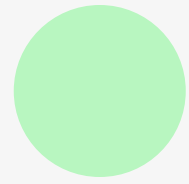
# GRU - D



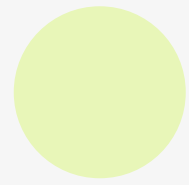
**Input Decay**

$$\gamma_t = \exp\{-\max(0, W_\gamma \delta_t + b_\gamma)\}$$

$$\hat{x}_t^d = m_t^d x_t^d + (1 - m_t^d) (\gamma_{x_t}^d x_{t'}^d + (1 - \gamma_{x_t}^d) \tilde{x}^d)$$

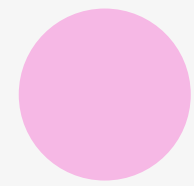


**Hidden State Decay**



**Mask Vector Inputs**

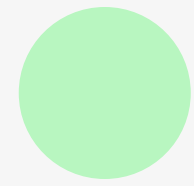
# GRU - D



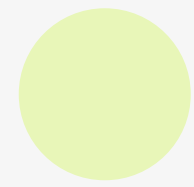
**Input Decay**

$$\gamma_t = \exp\{-\max(0, W_\gamma \delta_t + b_\gamma)\}$$

$$\hat{x}_t^d = m_t^d x_t^d + (1 - m_t^d) (\gamma_{x_t}^d x_{t'}^d + (1 - \gamma_{x_t}^d) \tilde{x}^d)$$



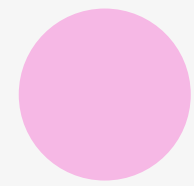
**Hidden State Decay**



**Mask Vector Inputs**



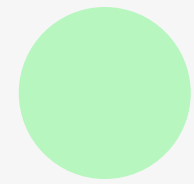
# GRU - D



## Input Decay

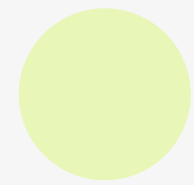
$$\gamma_t = \exp\{-\max(0, W_\gamma \delta_t + b_\gamma)\}$$

$$\hat{x}_t^d = m_t^d x_t^d + (1 - m_t^d) (\gamma_{x_t}^d x_{t'}^d + (1 - \gamma_{x_t}^d) \tilde{x}^d)$$



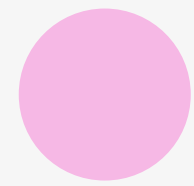
## Hidden State Decay

$$\hat{h}_{t-1} = \gamma_{h_t} \odot h_{t-1},$$



## Mask Vector Inputs

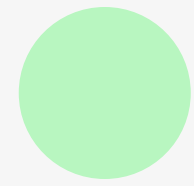
# GRU - D



**Input Decay**

$$\gamma_t = \exp\{-\max(0, W_\gamma \delta_t + b_\gamma)\}$$

$$\hat{x}_t^d = m_t^d x_t^d + (1 - m_t^d) (\gamma_{x_t}^d x_{t'}^d + (1 - \gamma_{x_t}^d) \tilde{x}^d)$$



**Hidden State Decay**

$$\hat{h}_{t-1} = \gamma_{h_t} \odot h_{t-1},$$



**Mask Vector Inputs**

**[0 0 0 1 1 1 0 0 0]**

## GRU Updates

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$$

$$\tilde{h}_t = \tanh(W x_t + U(r_t \odot h_{t-1}) + b)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

## GRU-D Updates

$$r_t = \sigma(W_r \hat{x}_t + U_r \hat{h}_{t-1} + V_r m_t + b_r)$$

$$z_t = \sigma(W_z \hat{x}_t + U_z \hat{h}_{t-1} + V_z m_t + b_z)$$

$$\tilde{h}_t = \tanh(W \hat{x}_t + U(r_t \odot \hat{h}_{t-1}) + V m_t + b)$$

$$h_t = (1 - z_t) \odot \hat{h}_{t-1} + z_t \odot \tilde{h}_t$$

## GRU Updates

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r)$$

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z)$$

$$\tilde{h}_t = \tanh(W x_t + U(r_t \odot h_{t-1}) + b)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$



**Input Decay**



**Hidden State Decay**



**Mask Vector Inputs**

## GRU-D Updates

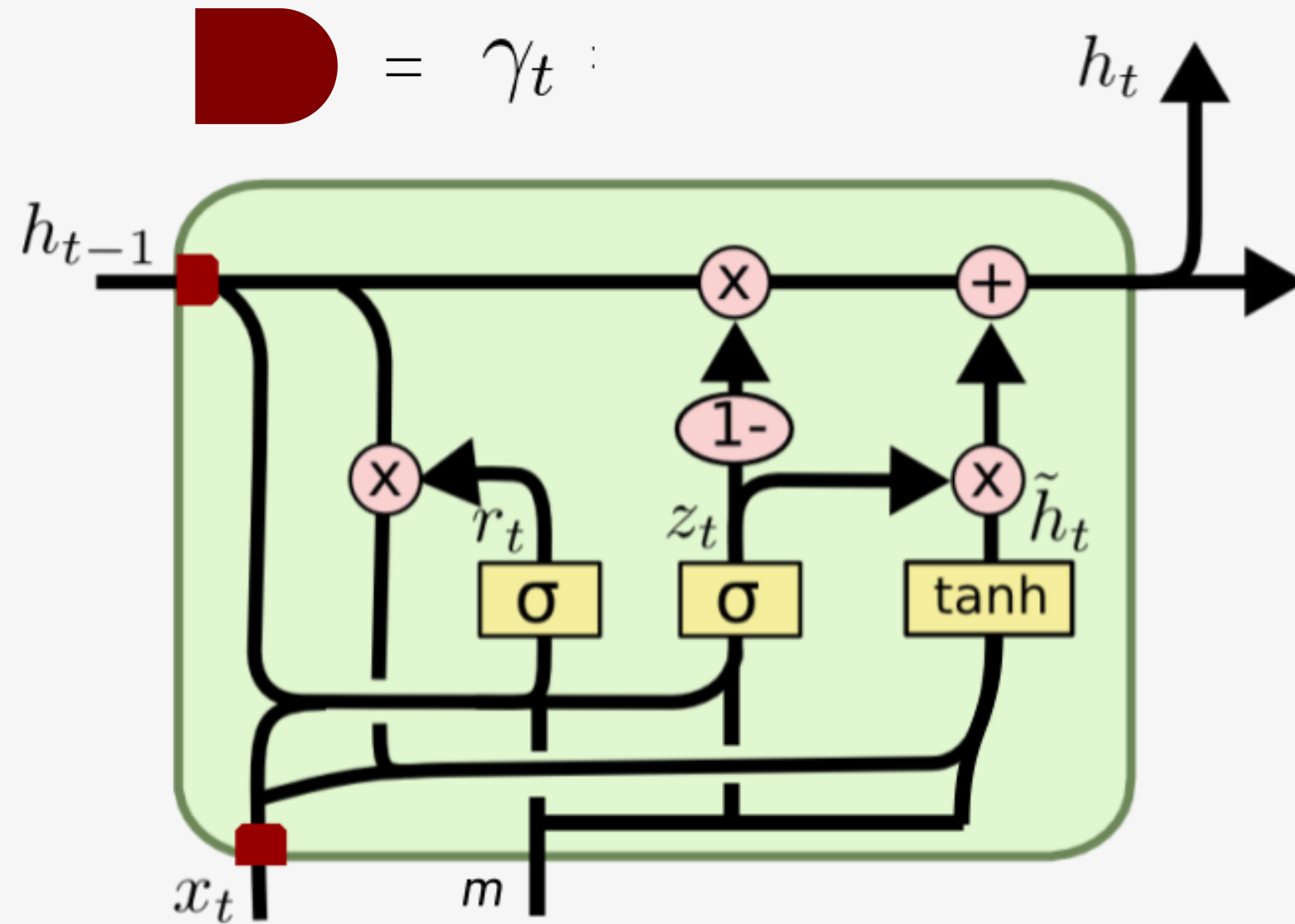
$$r_t = \sigma(\underbrace{W_r \hat{x}_t}_{\text{Input Decay}} + \underbrace{U_r \hat{h}_{t-1}}_{\text{Hidden State Decay}} + \underbrace{V_r m_t}_{\text{Mask Vector Inputs}} + b_r)$$

$$z_t = \sigma(\underbrace{W_z \hat{x}_t}_{\text{Input Decay}} + \underbrace{U_z \hat{h}_{t-1}}_{\text{Hidden State Decay}} + \underbrace{V_z m_t}_{\text{Mask Vector Inputs}} + b_z)$$

$$\tilde{h}_t = \tanh(\underbrace{W \hat{x}_t}_{\text{Input Decay}} + U(r_t \odot \underbrace{\hat{h}_{t-1}}_{\text{Hidden State Decay}}) + \underbrace{V m_t}_{\text{Mask Vector Inputs}} + b)$$

$$h_t = (1 - z_t) \odot \underbrace{\hat{h}_{t-1}}_{\text{Hidden State Decay}} + z_t \odot \tilde{h}_t$$

## GRU-D Updates



$$r_t = \sigma(\underbrace{W_r \hat{x}_t}_{\text{Input Decay}} + \underbrace{U_r \hat{h}_{t-1}}_{\text{Hidden State Decay}} + \underbrace{V_r m_t}_{\text{Mask Vector Inputs}} + b_r)$$

$$z_t = \sigma(\underbrace{W_z \hat{x}_t}_{\text{Input Decay}} + \underbrace{U_z \hat{h}_{t-1}}_{\text{Hidden State Decay}} + \underbrace{V_z m_t}_{\text{Mask Vector Inputs}} + b_z)$$

$$\tilde{h}_t = \tanh(\underbrace{W \hat{x}_t}_{\text{Input Decay}} + U(r_t \odot \underbrace{\hat{h}_{t-1}}_{\text{Hidden State Decay}}) + \underbrace{V m_t}_{\text{Mask Vector Inputs}} + b)$$

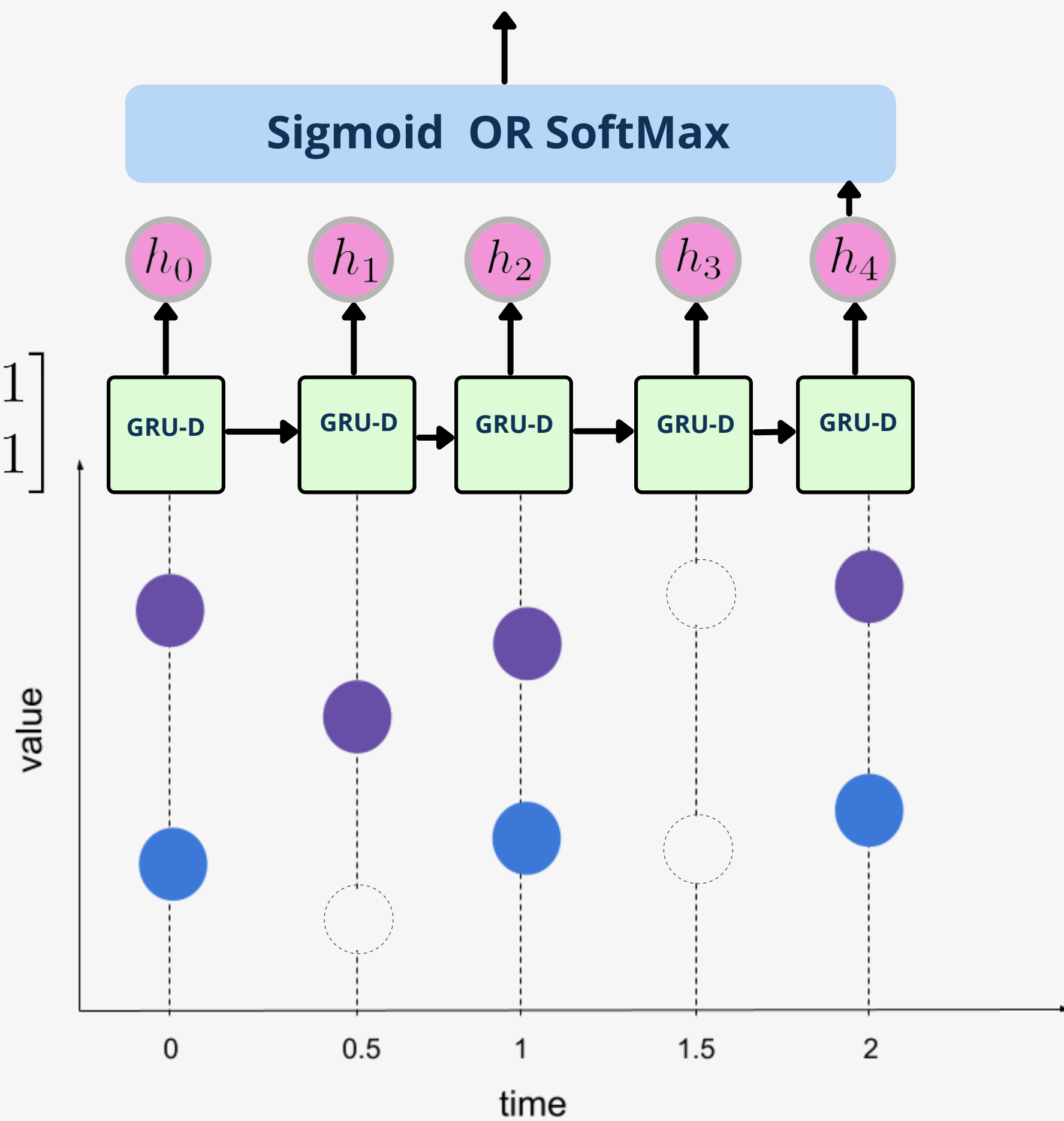
$$h_t = (1 - z_t) \odot \underbrace{\hat{h}_{t-1}}_{\text{Hidden State Decay}} + z_t \odot \tilde{h}_t$$

 **Input Decay**

 **Hidden State Decay**

 **Mask Vector Inputs**

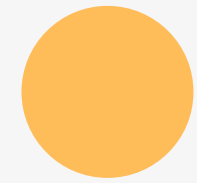
$$\mathbf{M} = \begin{bmatrix} 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \end{bmatrix}$$





## **Non-RNN Baselines**

Logistic Regression, SVM, Random Forest



## **GRU - Mean**

Missing values filled with feature mean



## **GRU - Forward**

Missing values filled with previous observed value



## **GRU - Simple**

Missingness masks and time since last observation passed as feature inputs



## **GRU - D**

Proposed model



## **Gesture Phase Segmentation**

Multivariate time series with synthetic missingness introduced - multiclass classification



## **PhysioNet Challenge 2012**

Multivariate time series ICU records

\*Binary classification (mortality)

\*Multiclass classification (4 tasks)



## **MIMIC III**

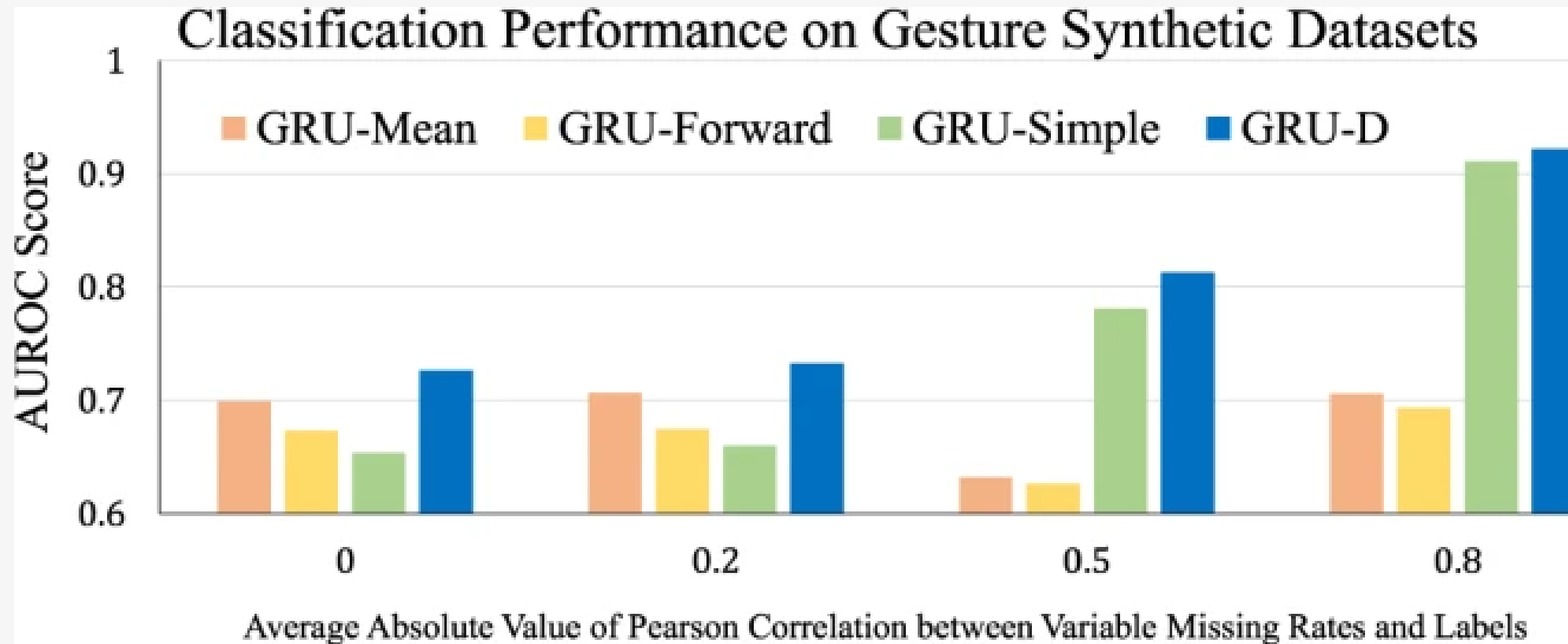
Multivariate time series ICU records

\*Binary classification (mortality)

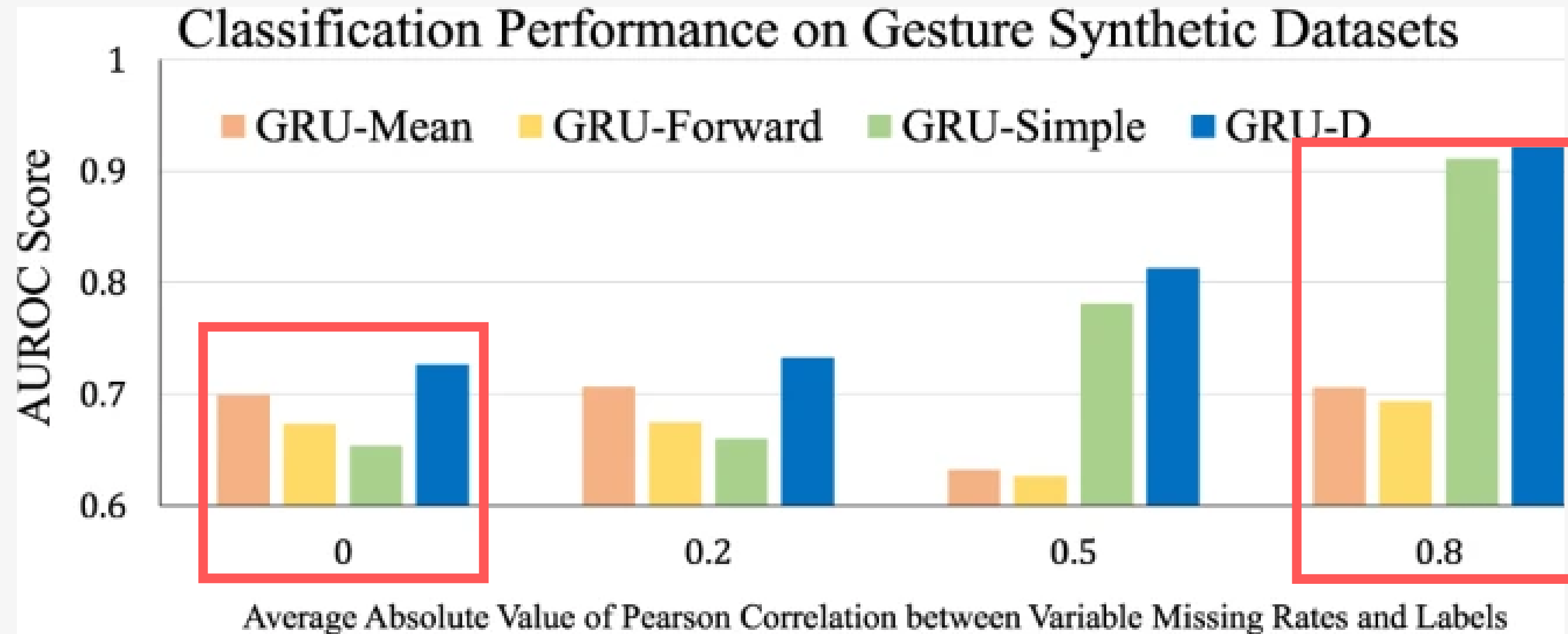
\*Multiclass classification (ICD9 codes, 20 tasks)



# Gesture Phase Segmentation



# Gesture Phase Segmentation



## PhysioNet (Mortality)

### MODEL

GRU - Mean  
GRU - Forward  
GRU - Simple  
★ GRU - D

### AUROC

$0.8252 \pm 0.011$   
 $0.8192 \pm 0.013$   
 $0.8380 \pm 0.008$   
 $0.8527 \pm 0.003$

## PhysioNet (4 tasks)

### MODEL

GRU - Mean  
GRU - Forward  
GRU - Simple  
★ GRU - D

### AUROC

$0.8099 \pm 0.011$   
 $0.8091 \pm 0.008$   
 $0.8249 \pm 0.010$   
 $0.8370 \pm 0.012$

## MIMIC (Mortality)

### MODEL

GRU - Mean  
GRU - Forward  
GRU - Simple  
★ GRU - D

### AUROC

$0.8162 \pm 0.014$   
 $0.8195 \pm 0.004$   
 $0.8226 \pm 0.010$   
 $0.8424 \pm 0.012$

## MIMIC (ICD 9 20 tasks)

### MODEL

GRU - Mean  
GRU - Forward  
GRU - Simple  
★ GRU - D

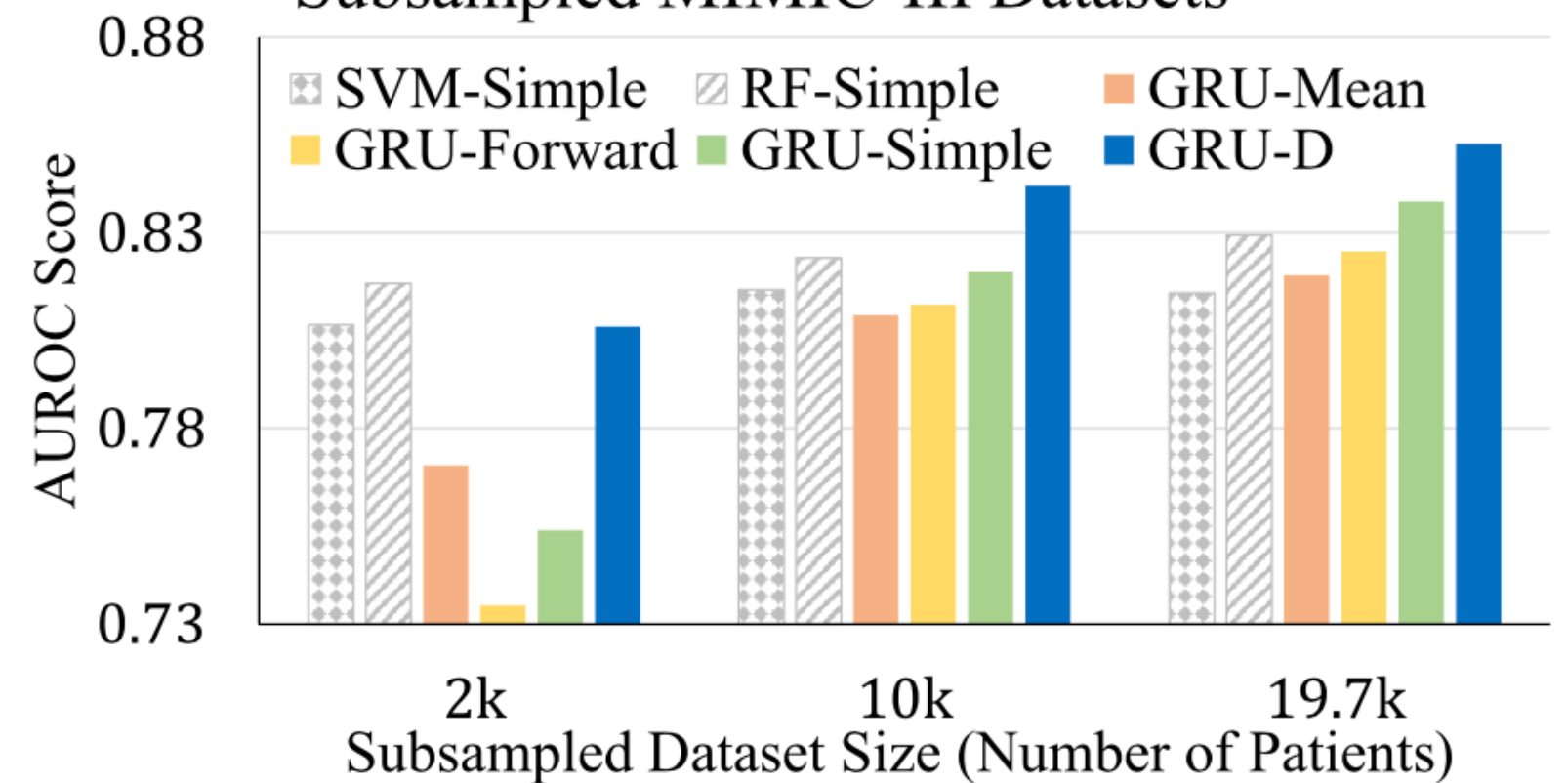
### AUROC

$0.7070 \pm 0.001$   
 $0.7077 \pm 0.001$   
 $0.7105 \pm 0.001$   
 $0.7123 \pm 0.003$

# Strengths

- Scalable
- Extensive evaluation
- Solution to data not missing-completely-at-random

Mortality Prediction Performance on Subsampled MIMIC-III Datasets



(b) Performance for predicting mortality on subsampled datasets.

# Limitations

- **Uninformative missingness**
  - No clear inherent correlation between the missingness pattern and prediction task.
- **Decay mechanism**
  - Needs to be explicitly designed for each domain.
  - Should it always decay?
- **Unsupervised settings**
  - Labels are required.

**Questions?**