

Motivation and key ideas

- The role of soft-DTW is to evaluate the (relaxed) DTW distance between a pair of sequences $\Psi \equiv [\psi_1, \dots, \psi_\tau] \in \mathbb{R}^{d' \times \tau}$, $\Psi' \equiv [\psi'_1, \dots, \psi'_{\tau'}] \in \mathbb{R}^{d' \times \tau'}$ of lengths τ and τ' , respectively. Under its transportation plan $\mathcal{A}_{\tau, \tau'}$, each path $\Pi \in \mathcal{A}_{\tau, \tau'}$ is evaluated to ascertain the path distance, and the smallest distance is 'selected' by the soft minimum: $d_{\text{DTW}}^2(\Psi, \Psi') = \text{SoftMin}_\gamma \left([\langle \Pi, D(\Psi, \Psi') \rangle]_{\Pi \in \mathcal{A}_{\tau, \tau'}} \right)$, where $\text{SoftMin}_\gamma(\alpha) = -\gamma \log \sum_i \exp(-\alpha_i/\gamma)$ is the soft minimum, $\gamma \geq 0$ controls its relaxation.
- However, the path distance $\langle \Pi, D(\Psi, \Psi') \rangle$ of path Π ignores the observation uncertainty of frame-wise feature representations by simply relying on the Euclidean distances stored in D .
- We model each path distance by the product of likelihoods of Normal distributions.

Our (soft) uncertainty-DTW takes the following generalized form:

$$\begin{cases} d_{\text{uDTW}}^2(D, \Sigma^\dagger) = \text{SoftMin}_\gamma \left([\langle \Pi, D \odot \Sigma^\dagger \rangle]_{\Pi \in \mathcal{A}_{\tau, \tau'}} \right) \\ \Omega(\Sigma) = \text{SoftMinSel}_\gamma \left(w, [\langle \Pi, \log \Sigma \rangle]_{\Pi \in \mathcal{A}_{\tau, \tau'}} \right), \end{cases} \quad (1)$$

$$\text{where } D \equiv D(\Psi, \Psi'), \Sigma \equiv \Sigma(\Psi, \Psi') \text{ and } \Sigma^\dagger = \text{inv}(\Sigma), \quad (2)$$

where \odot is the Hadamard product, $\Sigma^\dagger(\Psi, \Psi')$ is the element-wise inverse of matrix $\Sigma \in \mathbb{R}_+^{\tau \times \tau'} \equiv [\sigma^2(\psi_m, \psi'_n)]_{(m,n) \in \mathcal{I}_\tau \times \mathcal{I}_{\tau'}}$ which contains pair-wise variances between all possible pairings of frame-wise feature representations from sequences Ψ and Ψ' .

- Eq. (1) yields the uncertainty-weighted time warping distance $d_{\text{uDTW}}^2(D, \Sigma^\dagger)$ between sequences Ψ and Ψ' because D and Σ^\dagger are both functions of (Ψ, Ψ') .
- Eq. (2) provides the regularization penalty $\Omega(\Sigma)$ for sequences Ψ and Ψ' (as Σ is a function of (Ψ, Ψ')) which is the aggregation of log-variances along the path with the smallest distance, *i.e.*, path matrix $(\Pi_{i^*} \in \{0, 1\}^{\tau \times \tau'}: i^* = \arg \min_k w_k)$ if $\gamma = 0$, and vector w contains path-aggregated distances for all possible paths of the plan $\mathcal{A}_{\tau, \tau'}$.

- Below is an example of a generic similarity learning loss:

$$\arg \min_{\mathcal{P}} \sum \ell(d_{\text{uDTW}}^2(D(\Psi_n, \Psi'_n), \Sigma^\dagger(\Psi_n, \Psi'_n)), \delta_n) + \beta \Omega(\Sigma(\Psi_n, \Psi'_n)), \text{ or} \quad (3)$$

$$\arg \min_{\mathcal{P}, \Sigma > 0} \sum \ell(d_{\text{uDTW}}^2(D(\Psi_n, \Psi'_n), \Sigma^\dagger), \delta_n) + \beta \Omega(\Sigma), \quad (4)$$

where $\Psi_n = f(\mathbf{X}_n; \mathcal{P})$ and $\Psi'_n = f(\mathbf{X}'_n; \mathcal{P})$ are obtained from some backbone encoder $f(\cdot; \mathcal{P})$ with parameters \mathcal{P} and $(\mathbf{X}_n, \mathbf{X}'_n) \in \mathcal{X}$ is a sequence pair to compare with the similarity label $\delta_n \in \{0, 1\}$ (where $\delta_n = 0$ if $y_n = y'_n$ and $\delta_n = 1$ otherwise), (y_n, y'_n) is a pair of class labels for (Ψ_n, Ψ'_n) , and $\beta \geq 0$ controls the penalty for high matching uncertainty.

- Minimizing Eq. (4) w.r.t. (\mathcal{P}, Σ) assumes that $\Sigma \in \mathbb{R}_+^{\tau \times \tau'}$ is a free variable to minimize over. However, as sequence pairs vary in length, *i.e.*, $\tau \neq \tau'$, optimizing one global Σ is impossible (its size changes). Thus, we minimize loss functions with the distance/penalty in Eq. (3) and (4) where Σ is parametrized by (Ψ_n, Ψ'_n) :

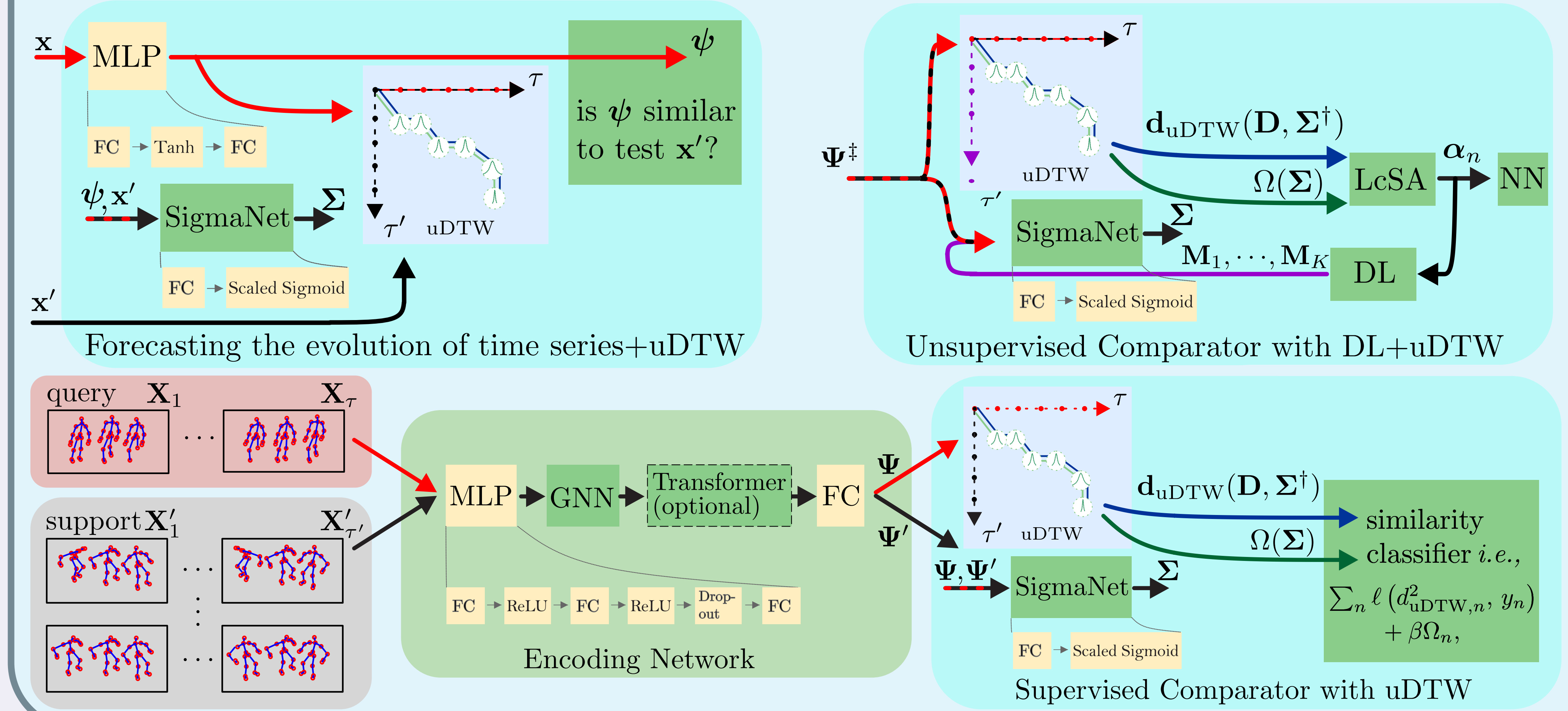
$$d_{\text{uDTW}}^2(\Psi, \Psi') \equiv d_{\text{uDTW}}^2(D(\Psi, \Psi'), \Sigma^\dagger(\Psi, \Psi')), \quad (5)$$

$$\Omega_\bullet(\Psi, \Psi') \equiv \Omega(\Sigma(\Psi, \Psi')). \quad (6)$$

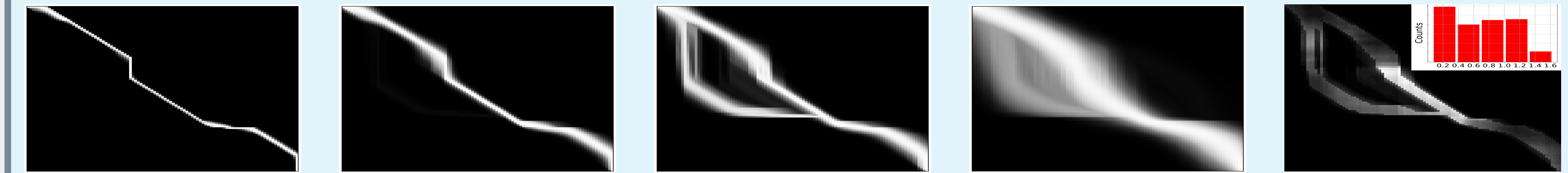
- We devise a small MLP (SigmaNet) for $\sigma(\cdot; \mathcal{P}_\sigma)$ or $\sigma(\cdot, \cdot; \mathcal{P}_\sigma)$ to obtain Σ (uses additive variance for individual frames ψ_m and ψ'_n) or Σ' (uses a jointly generated variance for (ψ_m, ψ'_n)).

The pipeline: further details

We define specific loss functions for several problems such as (i) forecasting the evolution of time series, (ii) clustering time series or (iii) even matching sequence pairs in few-shot action recognition.



Results



(a) $\text{sDTW}_{\gamma=0.01}$ (b) $\text{sDTW}_{\gamma=0.1}$ (c) $\text{uDTW}_{\gamma=0.01}$ (d) $\text{uDTW}_{\gamma=0.1}$ (e) uncertainty
We power-normalized pixels of plots (by 0.1) to see darker paths better. With higher γ that controls softness, in (b) & (d) more paths become 'active'. In (c), uDTW has two possible routes *vs.* sDTW (a) due to uncertainty modeling. In (e), we visualise uncertainty Σ . We binarize plot (c) and multiply it by the Σ to display uncertainty values on the path (white pixels = high uncertainty).

Few-shot action recognition results on NTU-60 & NTU-120:

#classes	10	20	30	40	50	#classes	20	40	60	80	100
Supervised						Supervised					
MatchNets	46.1	48.6	53.3	56.3	58.8	MatchNets	20.5	23.4	25.1	28.7	30.0
ProtoNet	47.2	51.1	54.3	58.9	63.0	ProtoNet	21.7	24.0	25.9	29.2	32.1
TAP	54.2	57.3	61.7	64.7	68.3	TAP	31.2	37.7	40.9	44.5	47.3
Euclidean	38.5	42.2	45.1	48.3	50.9	Euclidean	18.7	21.3	24.9	27.5	30.0
sDTW	53.7	56.2	60.0	63.9	67.8	sDTW	30.3	37.2	39.7	44.0	46.8
sDTW div.	54.0	57.3	62.1	65.7	69.0	sDTW div.	30.8	38.1	40.0	44.7	47.3
uDTW	56.9	61.2	64.8	68.3	72.4	uDTW	32.2	39.0	41.2	45.3	49.0
Unsupervised						Unsupervised					
Euclidean	20.9	23.7	26.3	30.0	33.1	Euclidean	13.5	16.3	20.0	24.9	26.2
sDTW	35.6	45.2	53.3	56.7	61.7	sDTW	20.1	25.3	32.0	36.9	40.9
sDTW div.	36.0	46.1	54.0	57.2	62.0	sDTW div.	20.8	26.0	33.2	37.5	42.3
uDTW	37.0	48.3	55.3	58.0	63.3	uDTW	22.7	28.3	35.9	39.4	44.0