

## Introduction

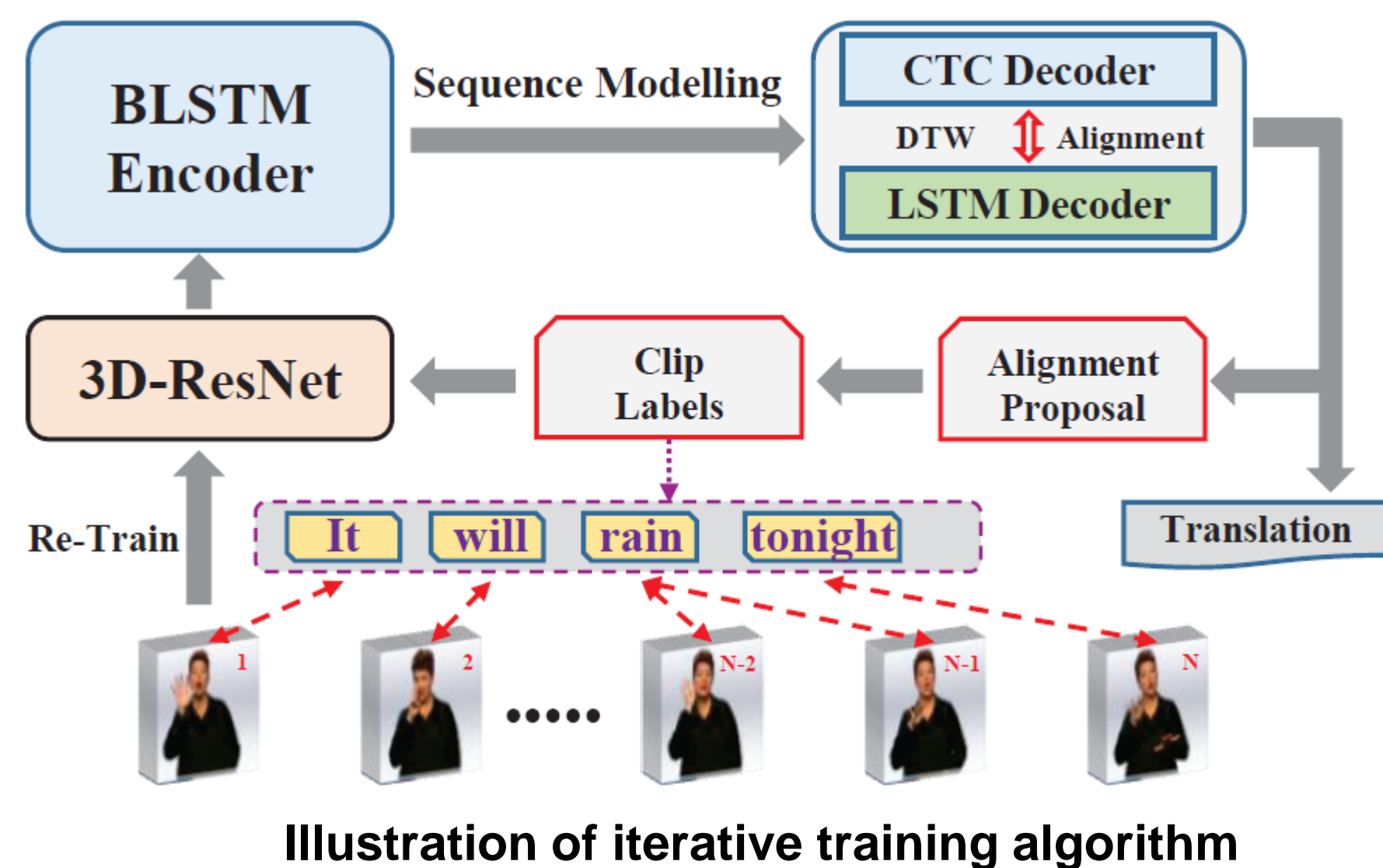
### Brief to SLR

Continuous sign language recognition (SLR) is a kind of weakly supervised sequence learning task, without rigid annotation of text words to video clips for a complete sign video. The key idea for continuous SLR is to learn the mapping between a sign video and its corresponding annotation of text sentence.

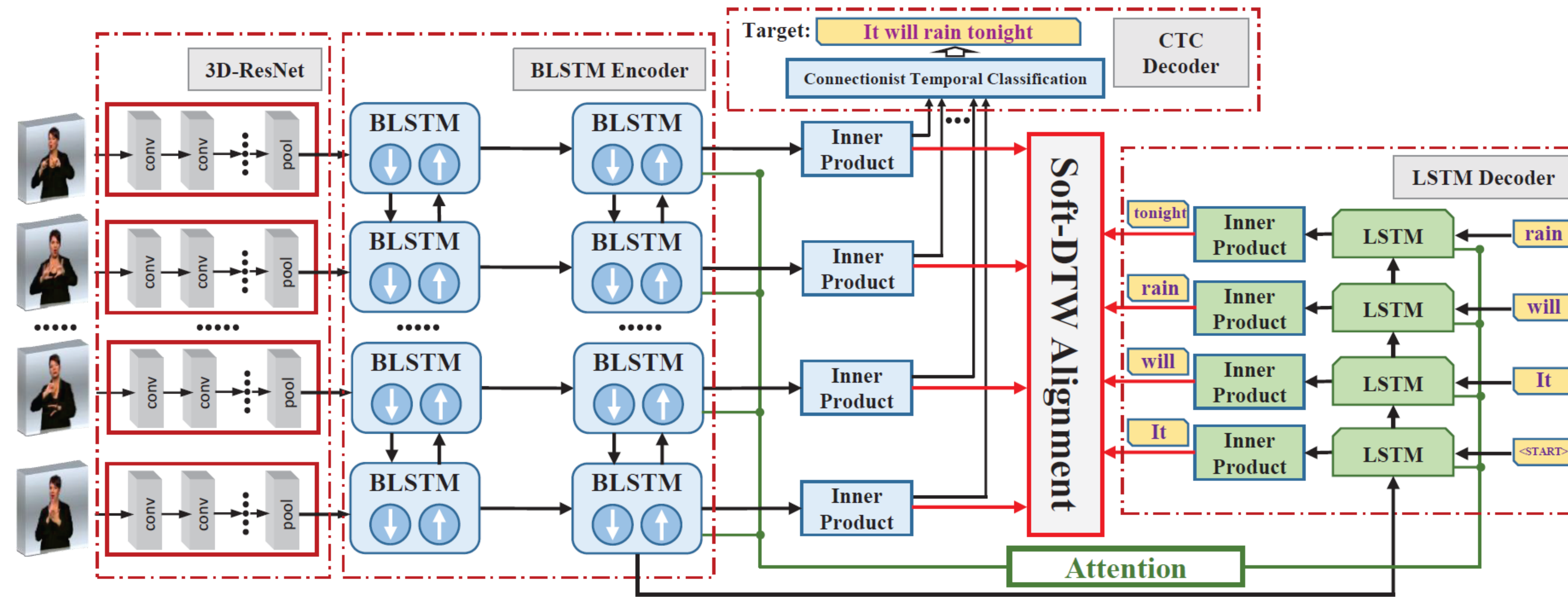
### Contribution

- A unified deep learning architecture integrating encoder-decoder network and connectionist temporal classification (CTC) for continuous SLR.
- A soft dynamic time warping (soft-DTW) alignment constraint between the LSTM CTC decoders, which indicates the temporal segmentation in sign videos.
- Iterative optimization strategy to train feature extractor and encoder-decoder network alternately with alignment proposals by warping path.

## Iterative Optimization



## Framework



### Modules

- Video Representation:  $\mathbf{F}^N = (f_1, \dots, f_N) = \{\mathcal{F}_\Theta(v_t)\}_{t=1}^N$
- Temporal Encoder:  $\mathbf{E} = \{e_t\}_{t=1}^N = \mathcal{R}(\{f_t\}_{t=1}^N)$
- LSTM Decoder:  $d_k = \text{Decoder}_{lstm}(c_k, s_k, h_{k-1}^d)$   
 $z_k = W_{fc2} \cdot d_k + b_{fc2} \quad \mathbf{Z} = (Z_{k,l}) = [z_1, z_2, \dots, z_L]^T$
- CTC Decoder:  $y_t = W_{fc1} \cdot e_t + b_{fc1} \quad \mathbf{Y} = (Y_{t,l}) = [y_1, y_2, \dots, y_N]^T$   
 $p(\pi|\mathbf{V}) = \prod_{t=1}^N p(\pi_t|v_t) = \prod_{t=1}^N Y_{t,\pi_t} \quad p_{ctc}(s|\mathbf{V}) = \sum_{\pi \in \mathcal{M}^{-1}(s)} p(\pi|\mathbf{V})$

- Representative features contribute to good performance.
- When training the network in an end-to-end way, the objective loss has limited contribution to the learning of parameters for low layers of feature extractor due to the chain rules of BP.
- We use soft-DTW alignment proposals as pseudo-labels to learn representative 3D-CNN features, and optimize feature extractor and sequence learning module in EM-like iterations.

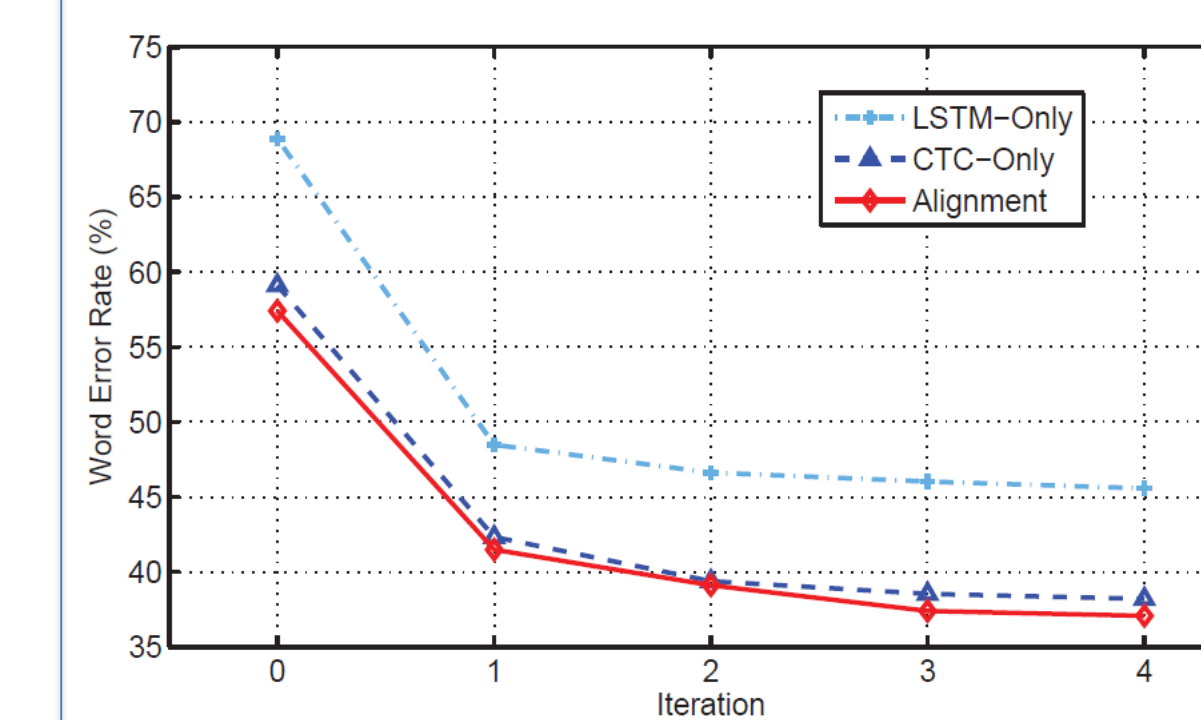
### Optimization and Decoding

- Soft-DTW:  $D_{i,j} = d_{i,j} + \min(D_{i-1,j}, D_{i,j-1}, D_{i-1,j-1})$   
 $d_{i,j} = \|u_i - v_j\|_2$   
 $\min_i \{a_i\} = \begin{cases} \min_i \{a_i\}, & \gamma = 0 \\ -\gamma \log \sum_i e^{-a_i/\gamma}, & \gamma > 0 \end{cases}$
- Objective Function:  
 $\mathcal{L} = \lambda \mathcal{L}_{ctc} + (1 - \lambda) \mathcal{L}_{lstm} + \mathcal{L}_{align} + \mu \|\omega\|^2$   
 $\mathcal{L}_{ctc} = -\ln p_{ctc}(s|\mathbf{V}) \quad \mathcal{L}_{lstm} = -\ln p_{lstm}(s|\mathbf{V}) \quad \mathcal{L}_{align} = \mathcal{D}_p(\mathbf{Y}, \mathbf{Z})$
- Decoding:  
 $r(s^i) = \alpha \ln p_{ctc}(s^i|\mathbf{V}) + (1 - \alpha) \ln p_{lstm}(s^i|\mathbf{V}) + \beta \ln L_i$

1. Extract features from 3D-ResNet for sign videos.
2. Train the encoder-decoder network with CTC by minimizing the total loss  $\mathcal{L}$ .
3. Get the warping path between the input clips and words provided by soft-DTW.
4. Fine-tune 3D-ResNet with the alignment as supervision.
5. Repeat 1-4. (Optimizing both modules in an EM-like way)

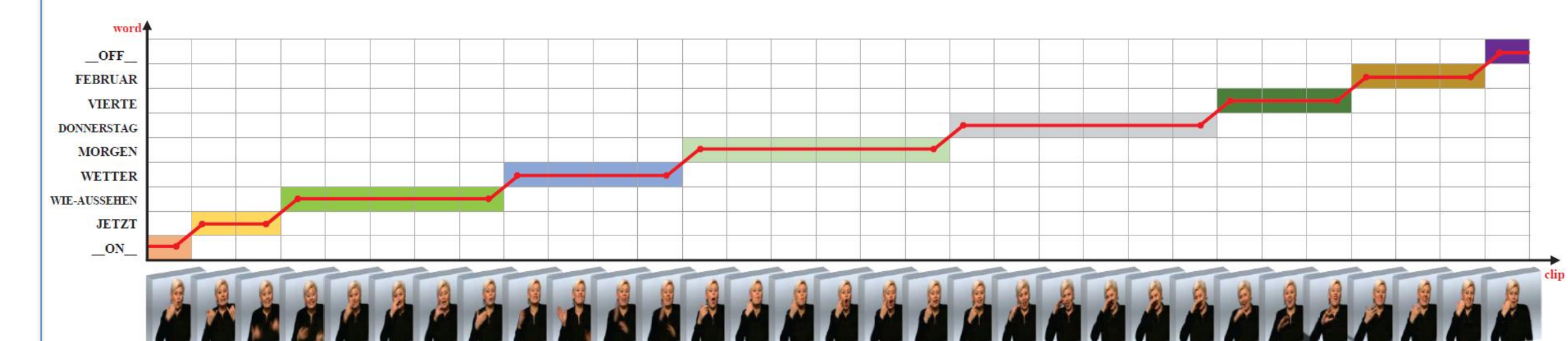
## Experiments

### □ RWTH-Phoenix-Weather (German Sign Language):



Performance comparison

### Alignment mechanism



### □ CSL-100 (Chinese Sign Language):

Method	Split I					Split II					
	Precision	BLEU-1	CIDEr	ROUGE-L	METEOR	Acc-w	BLEU-1	CIDEr	ROUGE-L	METEOR	WER
LSTM&CTC [12, 20]	0.858	0.936	8.632	0.940	0.646	0.332	0.343	0.241	0.362	0.111	0.757
S2VT [38]	0.897	0.902	8.512	0.904	0.642	0.457	0.466	0.479	0.461	0.189	0.670
S2VT (3-layer) [38]	0.903	0.911	8.592	0.911	0.648	0.461	0.475	0.477	0.465	0.186	0.652
HLSTM (SYS sampling) [15]	0.910	0.935	8.907	0.938	0.683	0.459	0.463	0.476	0.462	0.173	0.630
HLSTM [15]	0.924	0.942	9.019	0.944	0.699	0.482	0.487	0.561	0.481	0.193	0.662
HLSTM-attn [15]	0.929	0.948	9.084	0.951	0.703	0.506	0.508	0.605	0.503	0.205	0.641
Ours	0.939	0.980	9.342	0.981	0.713	0.670	0.724	3.946	0.716	0.383	0.327

### Performance comparison

## Conclusions

We propose a new deep architecture based on 3D-ResNet and encoder-decoder network with connectionist temporal classification by iterative optimization for continuous SLR.

- Jointly train encoder-decoder network by minimizing CTC loss and cross-entropy loss, with a soft-DTW alignment constraint.
- An iterative training strategy to optimize the 3D-ResNet and sequence learning module based on the warping path.
- Experimental results demonstrate the superiority of our method.