



PANIMALAR ENGINEERING COLLEGE



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

CS8811 PROJECT WORK

REVIEW NO:

Sign Language Recognition, Gloss Text, and Translation to Spoken Language: Utility for Aurally Challenged.

Guide Name: Dr. Senthil Kumar G
Batch No: E1

Team Members with Register number:
Isaac G (211418104091)
Jebastin Stephen Sarjin John (211418104096)

Project Details.

PROJECT TITLE	Sign Language Recognition, Gloss Text, and Translation to Spoken Language: Utility for Aurally Challenged.	
DOMAIN	NLP, Deep Learning, Neural Machine Transformer, Sign Language Translation.	
TEAM MEMBERS	Isaac G (211418104091)	Jebastin Stephen Sarjin John (211418104098)
PROJECT GUIDE	Dr. Senthil Kumar G	

Introduction.

Technology has aided today's communication very drastically. There are various means namely text, images, videos and other media by which information is shared and thereby communication is facilitated.

But a problem exists:

(1).Persons who are challenged with aural disorders and hearing disabilities find it difficult to understand audio information.

(2).Since there are diverse languages in the world, it would be difficult to translate one language to a language known to the particular person.

So in order to aid them, deep learning contributes features such as Continuous Sign Language Recognition (CSLR).

There would be some possible sign representations while communicating (i.e., hand shape, facial expression and body posture), Such signs has to be captured, analyzed and then converted to gloss language, which would be an intermediate between the Sign representation and the Spoken language.

We use two modules:

The Spatial-Modeling Cue (SMC) module is concerned with spatial representation and categorizes visual features of different cues with the aid of a self-contained pose estimation branch.

The Temporal Modeling Cue (TMC) module models temporal correlations along two parallel paths, i.e., intra-cue and inter- cue, which aims to preserve the uniqueness and explore the collaboration of multiple cues.

Finally, we design a joint optimization strategy to achieve the end-to-end sequence learning of the STMC network.

Literature Survey.

S.NO	Journal Type With Year	Authors	Title
1.	arXiv:2202.10419	Kayo Yin; Graham Neubig.	Interpreting Language Models with Contrastive Explanations.
2.	arXiv:2105.07476	Amit Moryossef; Kayo Yin; Graham Neubig; Yoav Goldberg.	Data Augmentation for Sign Language Gloss Translation.
3.	arXiv:2004.00588	Kayo Yin; Jesse Read.	Better Sign Language Translation with STMC-Transformer.
4.	arXiv:1409.0473	Bahdanau, D.; Cho, K.; and Bengio, Y. 2014.	Neural Machine Translation by jointly Learning to Align and Translate in ICLR.

Problem statement.

- This project is aimed for helping normal persons understand sign language with much ease.
- And also help the aurally challenged persons to understand the sign language.
- *eg:* Lets take a person A who is aurally challenged, to communicate he must either know sign language or he can't communicate. So to help the person A in both scenarios this project is developed. Through which sign language may be converted to spoken language to help the normal people understand, so that communication can take place.
- *eg:* Lets take person B who is normal person and he wants to either understand sign language or communicate to a deaf person, This project will help them understand what the aurally challenged person is saying.

Hardware/Software used.

SOFTWARE SPECIFICATION:

Technology Used: NLP, Deep Learning, OpenNMT (Neural Machine translation), CSLR (Continuous Sign Language Recognition).

Operating System: Windows 7 & above / Ubuntu

Software Requirement: Python 3.8, Anaconda

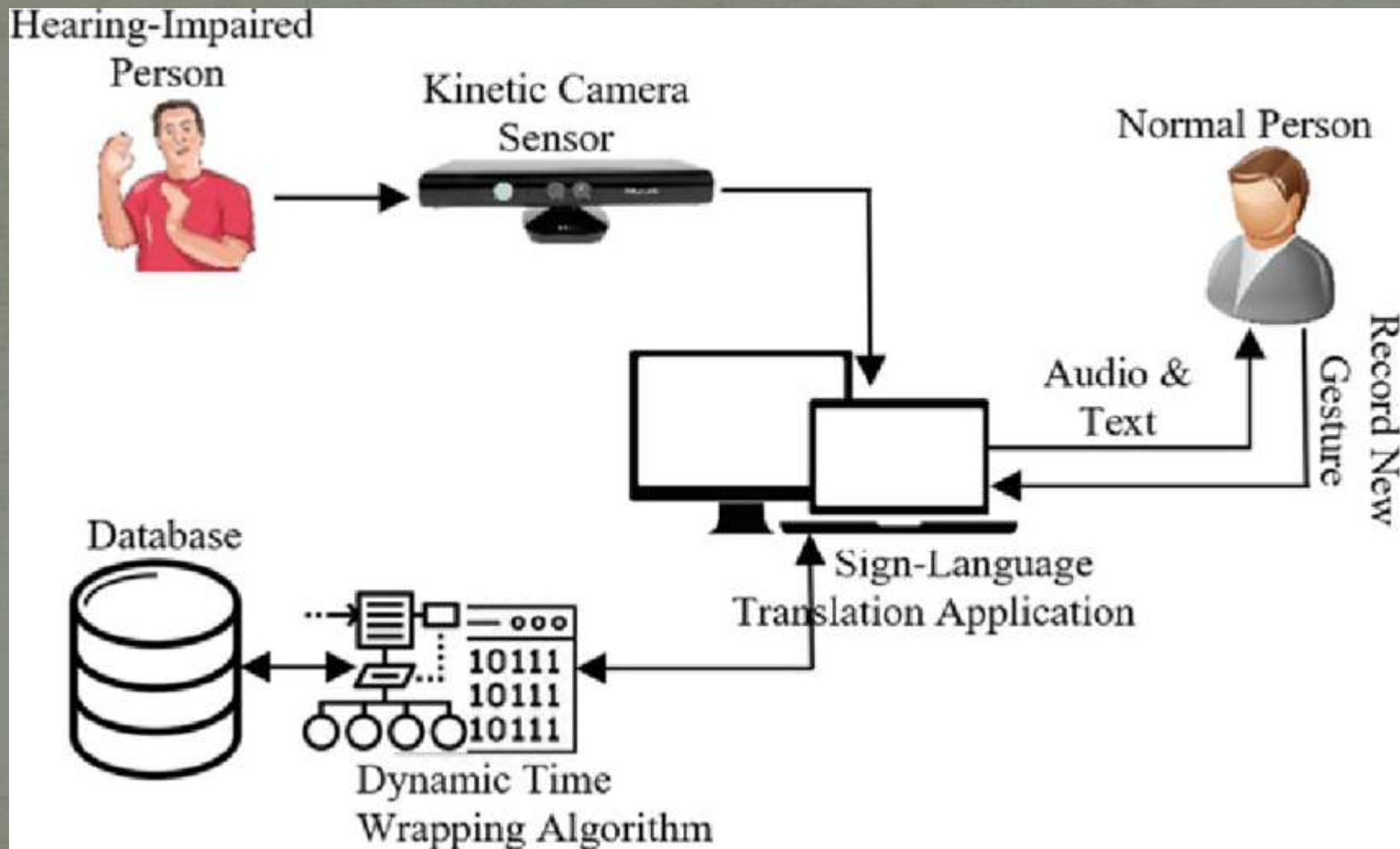
SYSTEM SPECIFICATIONS:

Ram: 8GB, 64-bit OS.

Processor: i3 6th Gen+ Intel Processor.

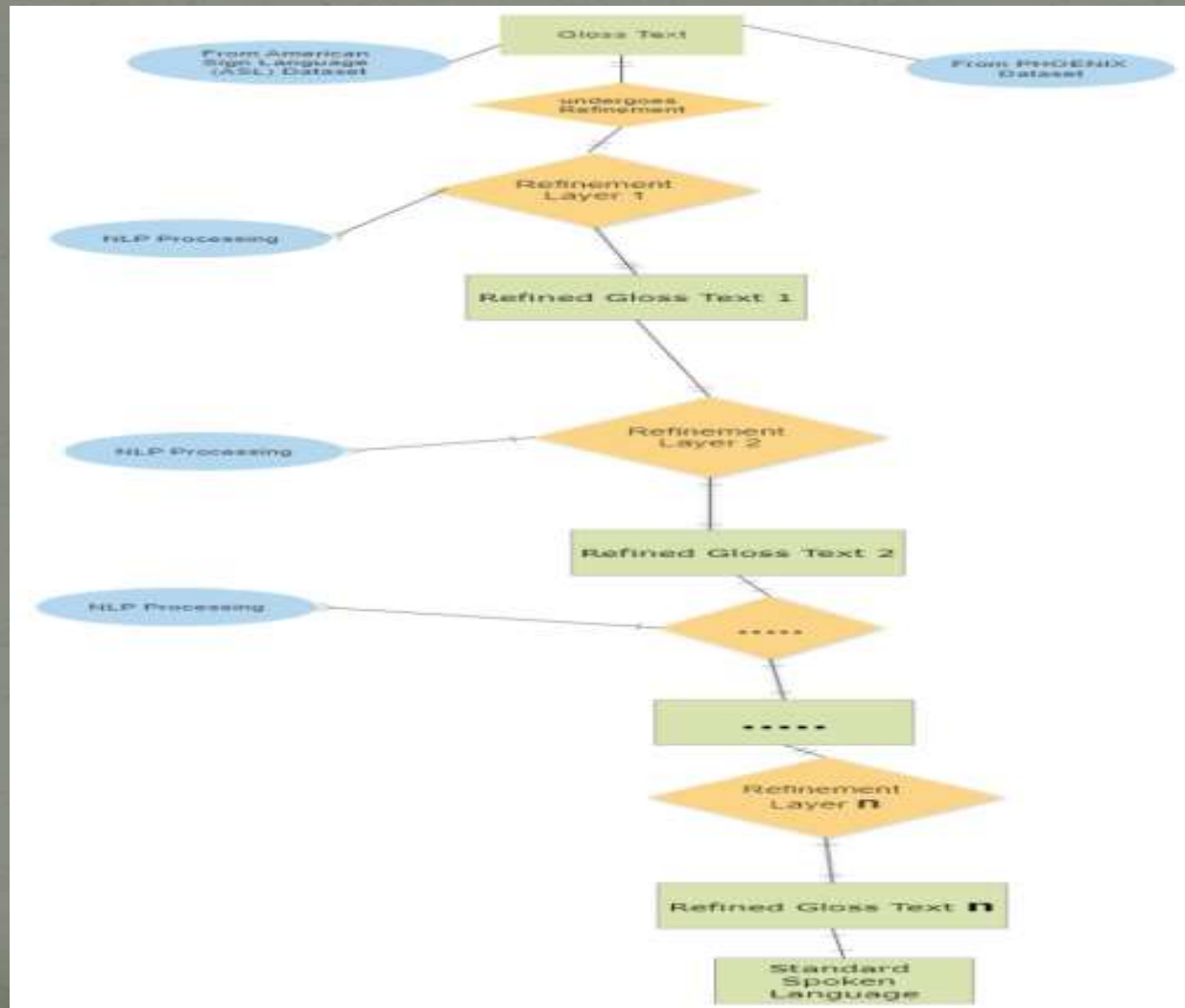
Hard Disk Capacity: 128 GB+

System Architecture.

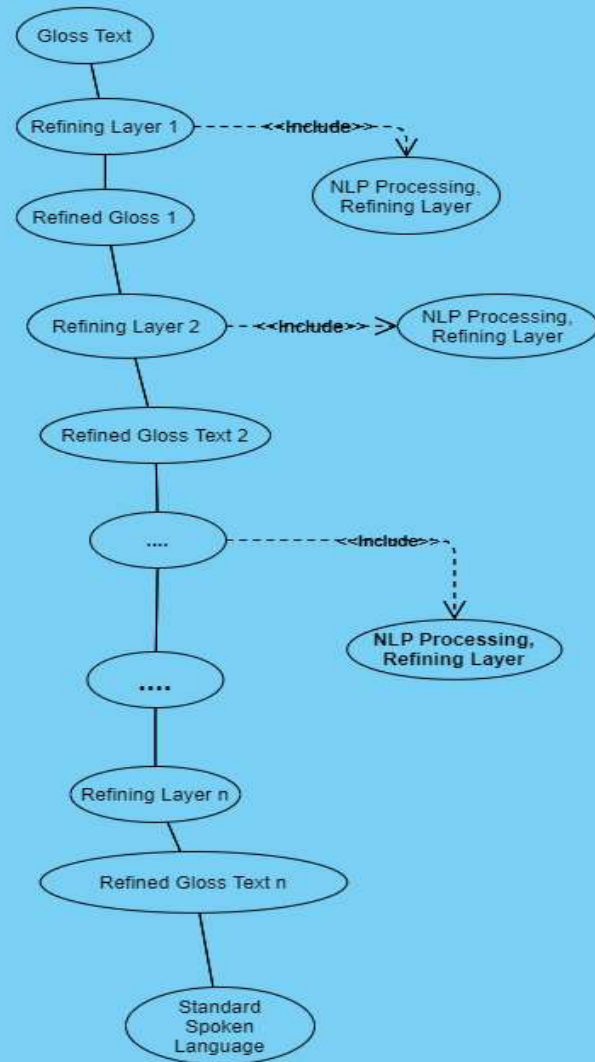


System Design.

(1.) ER Diagram:



(2.) Use-Case Diagram:



Modules.

(1.) Functionalities:

• **BLSTM Encoder:**

We use BLSTM Encoder to process the gloss text and the working of the encoder is given below;

With the help of BLSTM, The state transitions in the sequence of inputs has to be modeled using internal state of recurrent neural networks (RNN). Sign gloss sequence has to be mapped with spatial-temporal feature sequence using RNN.

RNN takes the feature sequence as input & then generates T^1 hidden state as:

$$\mathbf{h}_t = \mathbf{RNN}(\mathbf{h}_{t-1}, \mathbf{0}_t)$$

So, with the help of BLSTM Encoder, The Gloss text is processed.

• **Datasets:**

We Experiment the module with the help of 3 datasets:

1. PHOENIX-2014
2. CSL
3. PHOENIX-2014-T

The Data sets can be used to train, test, and validate the model.

(1.) PHOENIX-2014-T:

It is an extended version of PHOENIX-2014 and has two-stage annotations for new videos. One is sign gloss annotations for new videos. Another is German translation annotations for sign language translation (SLT) task. The videos are split into counts of 7096, 519 and 642 for Train, Dev and Test respectively. It has not overlap with the previous version between Train, Dev and Test set. The vocabulary size is 1115 for sign gloss and 300 for German.

(2.) CSL:

It is a Chinese Sign Language dataset, which has 100 sign language sentences about daily life with 178 words. Each sentence is performed by 50 signers and there are 5000 videos in total. For pre-training, it also provides a matched isolated Chinese Sign Language database, which contains 500 words. Each word is performed 10 times by 50 signers.

(3.) PHOENIX-2014:

It is a publicly available German Sign Language dataset, which is the most popular benchmark for CSLR. The corpus was recorded from broadcast new about the weather. It contains videos of 9 different signers with a vocabulary size of 1295. Videos are split into counts of 5672, 540 and 629 for Train, Dev and Test respectively. Our method is evaluated on the multi-signer database.

■ CSLR:

Computer Supported Learning Resource is function of performing continuous sign recognition in Real-time, Signs are sensed with CSLR.

In this project to perform experiments we use datasets of PHOENIX-2014, CSL, PHOENIX-2014-T.

CSLR is a task that recognizes the grammar of sign language itself, and it is difficult to provide meaningful interpretation only with CSLR because the grammar of sign language and the grammar of the spoken language are different.

■ Spatial multi-cue module:

The SMC module is dedicated to spatial representation and explicitly decomposes visual features of different cues with the aid of a self-contained pose estimation branch.

In SMC module we perform 3 actions: pose estimation, Patch cropping, Feature generation.

It consider the spatial movements or angles of a particular action.

■ Temporal Multi-cue model:

The Temporal Multi-cue (TMC) module models temporal correlations along two parallel paths, i.e., intra-cue and inter-cue, which aims to preserve the uniqueness and explore the collaboration of multiple cues. Finally, we design a joint optimization strategy to achieve the end-to-end sequence learning of the STMC network.

Temporal multi-cue (TMC) module intends to integrate spatiotemporal information from two aspects, intra-cue and inter-cue. The intra-cue path captures the unique features of each visual cue. The inter-cue path learns the combination of fused features from different cues at different time scales.

The TMC module contains two paths:

(1.) Intra-cue path:

Unique features of different cues at different time scales are provided. The temporal transformation inside each cue is performed as follows;

$$\begin{aligned} f_{l,n} &= \text{ReLU}(\mathcal{K}_k^{\frac{G}{N}}(f_{l-1,n})), \\ f_l &= [f_{l,1}, f_{l,2}, \dots, f_{l,N}]. \end{aligned}$$

(2.) Inter-cue path:

The temporal transformation has to be performed on the inter-cue feature from the previous block and fuse information from the intra-cue path as follows;

$$o_l = \text{ReLU}([\mathcal{K}_k^{\frac{Q}{2}}(o_{l-1}), \mathcal{K}_1^{\frac{Q}{2}}(f_l)]),$$

After each block, a temporal max-pooling with stride 2 and Kernel size 2 is performed. We use two blocks in the TMC module.

• Spatial Multi-cue module:

The SMC module is dedicated to spatial representation and explicitly decomposes visual features of different cues with the aid of a self-contained pose estimation branch.

It consider the spatial movements or angles of a particular action.

In SMC module we perform 3 actions:

- pose estimation
- Patch cropping
- Feature generation

(1.) Pose Estimation:

Deconvolutional networks are widely used in pixel-wise prediction. For pose estimation, Two deconvolutional layers of VGG-11. The stride of each layer is 2. So, the feature maps are $4x$ upsampled from the resolution $14x14$ to $56x56$. The output is fed into a point-wise convolutional layer to generate K predicted heat maps. In each heat map, the position of its corresponding keypoint is expected to show the highest response value.

(2.) Patch cropping:

In CSLR, the perception of detailed visual cues is vital, including eye gaze, facial expression, mouth shape, hand shape and orientations of hands. Our model takes predicted positions of the nose and both wrists as the center points of the face and both hands. The patches are cropped from the output of 4th convolutional layer of VGG-11. The cropping sizes are fixed to 24×24 for both hands and 16×16 for the face. It is large enough to cover body parts of a signer whose upper body is visible to the camera. The cropping sizes are fixed to 24×24 for both hands and 16×16 for the face. It is large enough to cover body parts of a signer whose upper body is visible to the camera. The center point of each patch is clamped into a range to ensure that the patch would not cross the border of the original feature map.

(3.) Feature Generation:

After K keypoints are predicted, They are flattened to a 1D-vector with dimension $2K$ and passed through two fully-connected (FC) layers with ReLU to get the feature vector of pose cue. Then, Feature maps of the face and both hands are cropped and processed by several convolutional layers, separately. Most sign gestures rely on the cooperation of the both hands. So we use weight-sharing convolutional layers for both hands. The outputs of them are concatenated along the channel-dimension. Finally, we perform global average pooling over all the feature maps with spatial dimension to form feature vectors of different cues.

Inference Time:

The inference time depends on the video length. In average, it takes around 8 seconds (25FPS) for a sign sentence. For a fair comparison, the inference time of 200 frames on a single GPU are evaluated. Compared with introducing an external VGG-11 based model for pose estimation, our self-contained branch saves around 44% inference time. It's notable that our framework with the self-contained branch still shows slightly better performance than an off-the-shelf model. Role regularization is facilitated by the differentiable pose estimation branch and thereby the overfitting of neural networks is alleviated.

Scoring:

We are using BLEU.py to find the success ratio or the speed at which the translation has been done. This metric will show how far is the efficiency of one work from another. Since this is a Research based project, we use this scoring system to prove that this project has better efficiency and accuracy compared to the previous model.

Screenshots Of Working.

(1.) Setup.py installation:

```
(one) D:\19-05\transformer-slt>python setup.py install
running install
C:\ProgramData\Anaconda3\envs\stone\lib\site-packages\setuptools\command\install.py:34: SetuptoolsDeprecationWarning: setup.py install is deprecated. Use build and pip and other standards-based tools.
  warnings.warn(
C:\ProgramData\Anaconda3\envs\stone\lib\site-packages\setuptools\command\easy_install.py:144: EasyInstallDeprecationWarning: easy_install command is deprecated. Use build and pip and other standards-based tools.
  warnings.warn(
running bdist_egg
running egg_info
creating OpenNMT_py.egg-info
writing OpenNMT_py.egg-info\PKG-INFO
writing dependency links to OpenNMT_py.egg-info\dependency_links.txt
writing entry points to OpenNMT_py.egg-info\entry_points.txt
writing requirements to OpenNMT_py.egg-info\requires.txt
writing top-level names to OpenNMT_py.egg-info\top_level.txt
writing manifest file 'OpenNMT_py.egg-info\SOURCES.txt'
reading manifest file 'OpenNMT_py.egg-info\SOURCES.txt'
adding license file 'LICENSE.md'
writing manifest file 'OpenNMT_py.egg-info\SOURCES.txt'
installing library code to build\bdist.win-amd64\egg
running install_lib
running build_py
creating build
creating build\lib
creating build\lib\onmt
copying onmt\model_builder.py -> build\lib\onmt
copying onmt\opts.py -> build\lib\onmt
copying onmt\trainer.py -> build\lib\onmt
copying onmt\train_single.py -> build\lib\onmt
copying onmt\_init_.py -> build\lib\onmt
creating build\lib\onmt\bin
copying onmt\bin\average_models.py -&> build\lib\onmt\bin
copying onmt\bin\preprocess.py -> build\lib\onmt\bin
copying onmt\bin\release_model.py -> build\lib\onmt\bin
copying onmt\bin\server.py -> build\lib\onmt\bin
copying onmt\bin\train.py -> build\lib\onmt\bin
copying onmt\bin\translate.py -> build\lib\onmt\bin
copying onmt\bin\_init_.py -> build\lib\onmt\bin
creating build\lib\onmt\decoders
copying onmt\decoders\decoder.py -> build\lib\onmt\decoders
copying onmt\decoders\ensemble.py -> build\lib\onmt\decoders
copying onmt\decoders\transformer.py -> build\lib\onmt\decoders
copying onmt\decoders\_init_.py -> build\lib\onmt\decoders
creating build\lib\onmt\encoders
copying onmt\encoders\encoder.py -> build\lib\onmt\encoders
copying onmt\encoders\mean_encoder.py -> build\lib\onmt\encoders
copying onmt\encoders\rnn_encoder.py -> build\lib\onmt\encoders
copying onmt\encoders\transformer.py -> build\lib\onmt\encoders
```

This is used to install necessary tools & packages for the application to process the Gloss text to Spoken language. The requirements will be installed in the virtual environment to perform the operation.

(2.) Data Processing:

```
(one) D:\19-05\transformer-ilt\onmt_preprocess -train_src data/phenix2014T.train.gloss -train_tgt data/phenix2014T.train.de -valid_src data/phenix2014T.dev.gloss -valid_tgt data/phenix2014T.dev.de -save_data data/dgs -lower
[2022-05-19 10:00:41,756 INFO] Extracting features...
[2022-05-19 10:00:41,771 INFO] * number of source features: 0.
[2022-05-19 10:00:41,771 INFO] * number of target features: 0.
[2022-05-19 10:00:41,771 INFO] Building 'Fields' object...
[2022-05-19 10:00:41,771 INFO] Building & saving training data...
[2022-05-19 10:00:43,053 INFO] * tgt vocab size: 2889.
[2022-05-19 10:00:43,053 INFO] * src vocab size: 1233.
[2022-05-19 10:00:43,053 INFO] Building & saving validation data...
```

It performs the task of converting data from a given form to a much more desired form. In this, the function of Training and Validity is performed.

(3.) Training Of The Model:

```
(one) D:\18-05\transformer-rlt\python train.py -data data/dgs -save_model model -keep_checkpoint 1 -layers 2 -rnn_size 512 -word_vec_size 512 -transformer_ff 2048 -heads 8 -encoder_type transformer -decoder_type transformer -position_encoding -max_generator_batches 2 -dropout 0.1 -early_stopping 1 -early_stopping_criteria accuracy ppl -batch_size 2048 -accum_count 1 -batch_type tokens -normalization tokens -optim adam -adam_beta2 0.998 -decay_method noam -warmup_steps 3600 -learning_rate 0.5 -max_grad_norm 0 -param_init 0 -param_init_glorot -label_smoothing 0.1 -valid_steps 100 -save_checkpoint_steps 100 -world_size 1
[2022-05-19 10:03:34,568 INFO] * src vocab size = 1233
[2022-05-19 10:03:34,568 INFO] * tgt vocab size = 2889
[2022-05-19 10:03:34,568 INFO] Building model...
[2022-05-19 10:03:34,740 INFO] NMTModel(
  (encoder): TransformerEncoder(
    (embeddings): Embedding(
      (make_embeddings): Sequential(
        (emb_luts): Elementwise(
          (0): Embedding(1233, 512, padding_idx=1)
        )
      )
    (pe): PositionalEncoding(
      (dropout): Dropout(p=0.1, inplace=False)
    )
  )
  (transformer): ModuleList(
    (0): TransformerEncoderLayer(
      (self_attn): MultiHeadedAttention(
        (linear_keys): Linear(in_features=512, out_features=512, bias=True)
        (linear_values): Linear(in_features=512, out_features=512, bias=True)
        (linear_query): Linear(in_features=512, out_features=512, bias=True)
        (softmax): Softmax(dim=-1)
        (dropout): Dropout(p=0.1, inplace=False)
        (final_linear): Linear(in_features=512, out_features=512, bias=True)
      )
      (feed_forward): PositionwiseFeedForward(
        (w_1): Linear(in_features=512, out_features=2048, bias=True)
        (w_2): Linear(in_features=2048, out_features=512, bias=True)
        (layer_norm): LayerNorm((512,)), eps=1e-06, elementwise_affine=True
        (dropout_1): Dropout(p=0.1, inplace=False)
        (relu): ReLU()
        (dropout_2): Dropout(p=0.1, inplace=False)
      )
      (layer_norm): LayerNorm((512,)), eps=1e-06, elementwise_affine=True
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (1): TransformerEncoderLayer(
      (self_attn): MultiHeadedAttention(
        (linear_keys): Linear(in_features=512, out_features=512, bias=True)
        (linear_values): Linear(in_features=512, out_features=512, bias=True)
        (linear_query): Linear(in_features=512, out_features=512, bias=True)
        (softmax): Softmax(dim=-1)
        (dropout): Dropout(p=0.1, inplace=False)
        (final_linear): Linear(in_features=512, out_features=512, bias=True)
      )
      (feed_forward): PositionwiseFeedForward(
        (w_1): Linear(in_features=512, out_features=2048, bias=True)
        (w_2): Linear(in_features=2048, out_features=512, bias=True)
        (layer_norm): LayerNorm((512,)), eps=1e-06, elementwise_affine=True
        (dropout_1): Dropout(p=0.1, inplace=False)
        (relu): ReLU()
        (dropout_2): Dropout(p=0.1, inplace=False)
      )
      (layer_norm): LayerNorm((512,)), eps=1e-06, elementwise_affine=True
      (dropout): Dropout(p=0.1, inplace=False)
    )
  )
)
```

This shows the models through which the words are trained for conversion of sign Language by using NLP(Natural Language Processing) & Hyper-parameter Tuning.

(4.) Output of the Trained model:

```
2022-05-19 11:01:15,690 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:15,690 INFO number of samples: 7000
2022-05-19 11:01:15,756 INFO Step 1150/100000; acc: 77.49; ppl: 1.40; mem: 1.46; lr: 0.00016; 1104/1044 tok/s; 3220 sec
2022-05-19 11:01:15,756 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:15,776 INFO number of samples: 7000
2022-05-19 11:01:15,872 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:15,872 INFO number of samples: 7000
2022-05-19 11:01:16,161 INFO Step 1400/100000; acc: 75.93; ppl: 1.60; mem: 0.94; lr: 0.00016; 1105/1039 tok/s; 3053 sec
2022-05-19 11:01:16,162 INFO loading dataset from data\ds\valid.0.pt
2022-05-19 11:01:16,177 INFO number of samples: 519
2022-05-19 11:01:16,284 INFO validation perplexity: 10.6255
2022-05-19 11:01:16,344 INFO validation accuracy: 50.0321
2022-05-19 11:01:16,344 INFO Stalled patience: 2/3
2022-05-19 11:01:16,360 INFO Saving checkpoint model_step_1400.pt
2022-05-19 11:01:16,374 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:16,374 INFO number of samples: 7000
2022-05-19 11:01:16,463 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:16,463 INFO number of samples: 7000
2022-05-19 11:01:16,574 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:16,574 INFO number of samples: 7000
2022-05-19 11:01:16,672 INFO Step 1650/100000; acc: 77.67; ppl: 1.44; mem: 0.84; lr: 0.00009; 1106/1044 tok/s; 3107 sec
2022-05-19 11:01:16,672 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:16,772 INFO number of samples: 7000
2022-05-19 11:01:16,872 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:16,872 INFO number of samples: 7000
2022-05-19 11:01:16,972 INFO Step 1900/100000; acc: 79.48; ppl: 1.28; mem: 0.81; lr: 0.00009; 1110/1044 tok/s; 3010 sec
2022-05-19 11:01:16,972 INFO loading dataset from data\ds\valid.0.pt
2022-05-19 11:01:16,972 INFO number of samples: 519
2022-05-19 11:01:16,999 INFO validation perplexity: 10.0542
2022-05-19 11:01:17,000 INFO validation accuracy: 50.0364
2022-05-19 11:01:17,000 INFO Stalled patience: 1/3
2022-05-19 11:01:17,015 INFO Saving checkpoint model_step_1900.pt
2022-05-19 11:01:17,015 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:17,081 INFO number of samples: 7000
2022-05-19 11:01:17,103 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:17,103 INFO number of samples: 7000
2022-05-19 11:01:17,270 INFO Step 2150/100000; acc: 81.27; ppl: 1.14; mem: 0.76; lr: 0.00011; 1119/1075 tok/s; 4050 sec
2022-05-19 11:01:17,270 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:17,368 INFO number of samples: 7000
2022-05-19 11:01:17,468 INFO loading dataset from data\ds\train.0.pt
2022-05-19 11:01:17,468 INFO number of samples: 7000
2022-05-19 11:01:17,568 INFO Step 2400/100000; acc: 83.32; ppl: 1.00; mem: 0.66; lr: 0.00022; 1130/1007 tok/s; 4170 sec
2022-05-19 11:01:17,568 INFO loading dataset from data\ds\valid.0.pt
2022-05-19 11:01:17,568 INFO number of samples: 519
2022-05-19 11:01:17,604 INFO validation perplexity: 11.4549
2022-05-19 11:01:17,604 INFO validation accuracy: 50.1339
2022-05-19 11:01:17,604 INFO Stalled patience: 0/3
2022-05-19 11:01:17,640 INFO Training finished after stalled validations, Early Stop!
2022-05-19 11:01:17,640 INFO Best model saved at step 1200
2022-05-19 11:01:17,640 INFO Saving checkpoint model_step_1900.pt
(www) D:\V9-09\transformer-gpt-
```

This shows that the training of the model has been completed and that the words have been learned.

Conclusion.

- Finally, We present a novel multi-cue framework for CSLR, which aims to learn spatial-temporal correlations of visual cues in an end-to-end fashion.
- In our framework, a spatial multi-cue module is designed with a self-contained pose estimation branch to decompose spatial multi-cue features.
- Moreover, we propose a temporal multi-cue module composed of the intra-cue and inter-cue paths, which aims to preserve the uniqueness of each cue and explore the synergy of different cues at the same time.
- A joint optimization strategy is proposed to accomplish multi-cue sequence learning.
- Extensive experiments on three large-scale CSLR datasets demonstrate the superiority of our STMC framework.

References:

- Kayo Yin; Graham Neubig. **Interpreting Language Models with Contrastive Explanations.**
- Amit Moryossef; Kayo Yin; Graham Neubig; Yoav Goldberg. **Data Augmentation For Sign Language Gloss Translation.**
- Kayo Yin; Jesse Read. **Better Sign Language Translation With STMC-Transformer.**
- Bahanau, D.; Cho, K.; and Bengio, Y. 2014. **Neural machine Translation by Jointly Learning to align and translate in ICLR.**
- Buehler, P.; Zisserman, A.; and Everingham, M. 2009. **Learning Sign Language by watching tv (using weakly aligned subtitles).** In *CVPR*.
- Chapelle, O., and Wu, M. 2010. **Gradient descent optimization of smoothed information retrieval metrics.** *Information Retrieval*.
- CihanCamgoz, N.; Hadfield, S.; Koller, o.; and Bowden, R. 2017. **Subunets: end-to-end hand shape and continuous sign language recognition.** In *ICCV*.
- CihanCamgoz, N.; Hadfield, s.; Koller, o.; Ney, H.; and Bowden, R. 2018. **Neural sign language translation.** In *CVPR*.
- Cooper, H., and Bowden, R. 2009. **Learning signs from subtitles: A weakly supervised approach to sign language recognition.** In *CVPR*
- Cui, R.; Liu, H.; and Zhang, c. 2017. **Recurrent convolutional neural networks for continuous sign language recognition byu staged optimization.** In *CVPR*.
- Cui, R.; Liu, H.; and Zhang, C. 2019. **A deep neural framework for continuous sign language recognition by iterative training.** *TMM* 21(7):1880-1891.
- Feichtenhofer, c.; Pinz, A.; and Zisserman, A. 2016. **Convolutional two-stream network fusion for video action recognition.** In *CVPR*.