

Lab 4 - Conditional VAE for Video Prediction

Cheng-yuan Ho

- Assignment release: 2024/7/30 18:30
- Assignment announce: 2024/7/30 18:30
- New E3 Deadline: **2024/8/13 18:00**
- Kaggle Deadline: **2024/8/13 11:55**
- Demo: 2024/8/13 **after Lab5 announcement (about 19:00)**
- Format:
 - Zip whole source code directory and named it in **LAB4_{studentID}_{YOUR_NAME}.zip**
 - Save your report as pdf file and named it in **LAB4_{studentID}_{YOUR_NAME}.pdf**
- About kaggle
 - Team name: {student id}_{your name}
 - 1 person 1 team
 - 5 submission per day

1. Introduction

In this assignment, we need to implement conditional video prediction in a VAE-based model. Before we go through this LAB, I want to mention that the topic about this LAB is highly correlated to the ICCV paper [1] which makes the prediction in a GAN-based model and ICML paper [2] which makes the video prediction in the VAE-based model. It might be helpful to study these papers [1], [2] before doing your work.

In [1], the author posts an interesting approach that predicts the video clips in GAN model structure. By taking a pose image generated by a pre-trained pose estimation network and a previous frame, the model can generate the next consecutive frame with a comparable subjective quality. Further details can be found in [1].



Figure 1. Taking the pose image (left down) that is generated by an off-the-shelf pose estimation network as an input to generate the next consecutive frame (right image).

In [2] , the author used a VAE-based model which combined the LSTM module to predict future frames in RNN manner. By taking two reference frames, the model has the ability to predict the following future frames. If you have NO ideas about how to perform video generation, reading this paper [2] will provide you great insights.

2. Lab Description

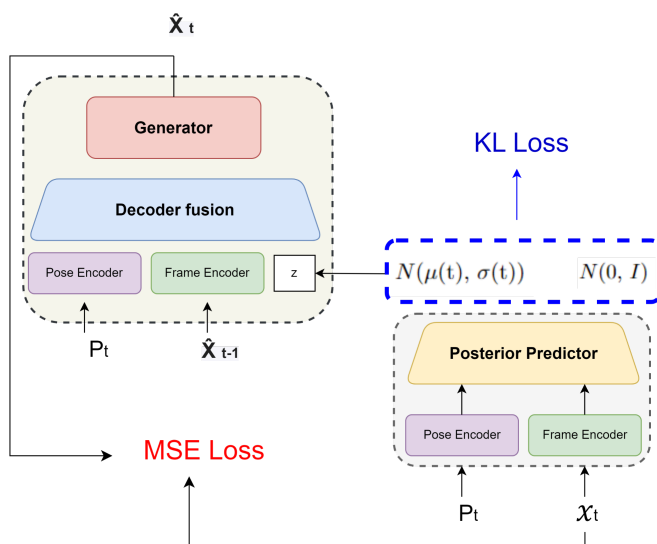


Figure 2. (a)

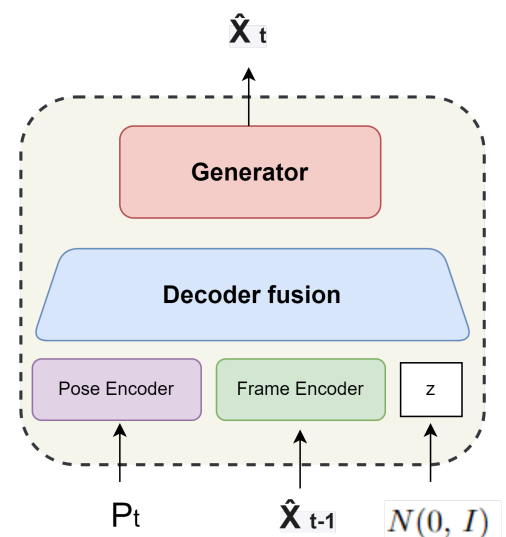


Figure 2. (b)

Figure 2. (a) give you a coarse version of the training protocol. Posterior Predictor which takes the current frame and current label as input to generate the distribution. The generator which takes the current label, last generated frame, noise sampled by the distribution predicted by Posterior Predictor as inputs to generate the current frame. Figure 2. (b) is the video generation protocol in inference time that the noise will directly sample from the prior distribution.

3. Requirements

a. Training details implementation

- i. Implement Video prediction protocol in training/testing stage
- ii. Implement reparameterization tricks
- iii. KL annealing implementation. (a) Cyclical. (b) Monotonic

b. Teacher forcing strategy

- i. Setup teacher forcing ratio, and plot the diagram in training
- ii. Teacher forcing ratio: 0~1

c. Other training strategy (training trick)

d. Plot diagrams

- i. Plot the loss curve while training
- ii. Plot PSNR-per frame diagram while validation

e. Analysis

- i. Compare the result in the loss curve if you apply different KL annealing strategies or even without KL annealing. Which one is better ?

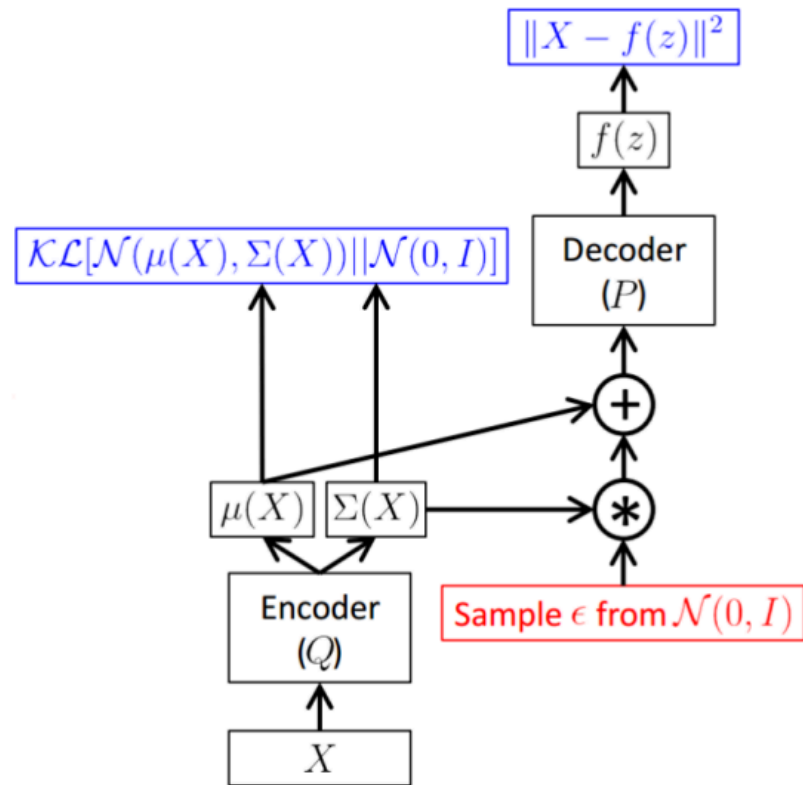
- f. Validate your result and make it into gif file
 - i. Pick one video clips in testing dataset and make the frames generated by your model into a gif file (This should be shown to TAs in LAB4 demo)
- g. Training strategy (Other training tricks)
 - i. Describe your training strategy in detail
- h. Derivation of Conditional VAE
 - i. Hand write the derivation of conditional VAE in your report
The objective function of conditional VAE is formulated as the following. You can refer to the EM algorithm from L13 slide 23.

$$L(X, c, q, \theta) = E_{z \sim q(Z|X, c; \phi)} \log p(X|Z, c; \theta) - KL(q(Z|X, c; \phi) || p(Z|c)) \quad (2)$$

4. Implementation Details

a. VAE recap

In the DL lecture, VAE has been explained thoroughly. While training VAE N2N, we need to adopt a reparameterization trick. For the purpose of stable training, the output of the reparameterization trick should be **log variance** instead of variance directly.



$$\underbrace{E_{\mathbf{Z} \sim q(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}')} \log p(\mathbf{X}|\mathbf{Z}; \boldsymbol{\theta})}_{\text{Re-parameterization for end-to-end training}} - \text{KL}(q(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}') || p(\mathbf{Z}))$$

Figure 2.1 The illustration of the reparameterization trick.

- a. **In one training step:** There will be k (default=16) frames as a training video sequence $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$, and 16 label frames as conditional signals $\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k\}$. First frame is the past frame which is provided to predict the consecutive $k-1$ future frames. Example is provided in Figure 3.

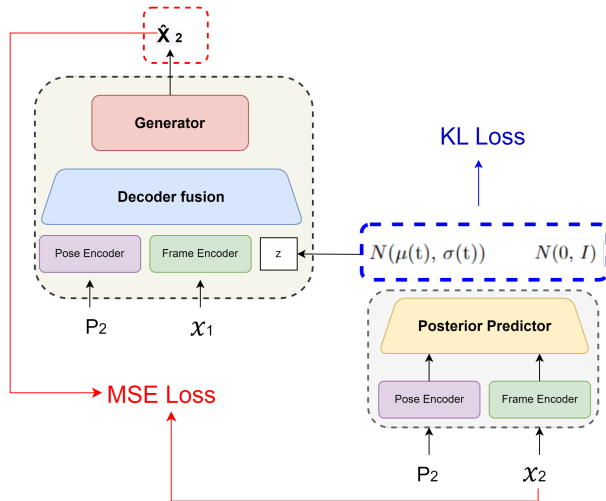


Figure 3. (a)

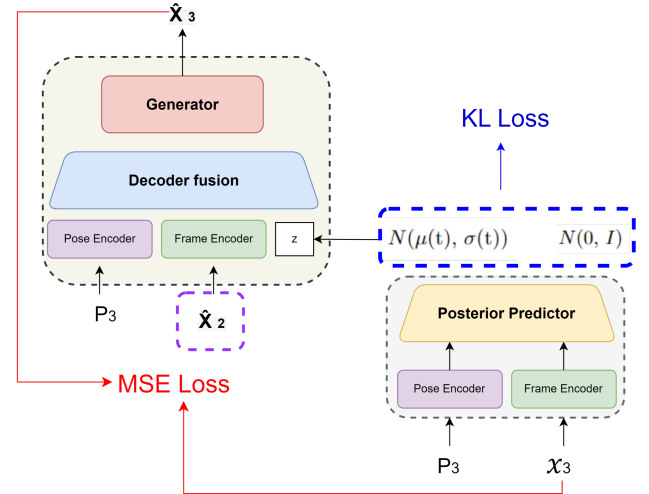


Figure 3. (b)

Figure 3. (a), and (b) give you some examples of how to generate the future frames. In (a), the task is to generate frame X_2 . However, we have no last frame to be taken as Decoder fusion input in scenario (a). Hence, X_1 will be provided in the dataset to be the first input of the generative system. In (b), the task is to generate x_3 . Decoder fusion takes the frame generated by the last step (red dash box in fig. 3(a)) as input to make the prediction.

- b. **In inference time:** $\{ X_1, P_2, Z \}$ will first be taken as input to generate the second frame. After the second frame is generated, it will be taken as an input to generate the next consecutive frame. Further detail can be found in Figure 4.

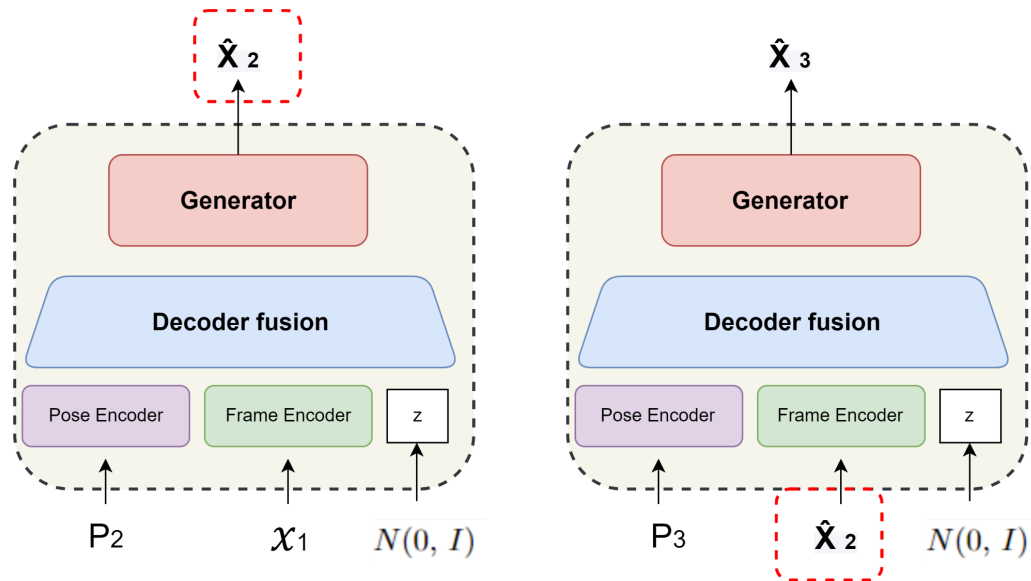


Figure 4 Example of how inference works

There would be 630 consecutive frames in the validation dataset. By taking one past frame, your model should predict **the remaining 629 frames**. After the prediction, it is suggested that you convert these 629 frames into a gif file, and see the quality of your prediction.

- c. **KL annealing strategy:** A variable weight is added to the KL term in the loss function. The following experiment results should be conducted and make the comparison in your report. It is suggested that you plot the loss curve in training while applying different strategies. Further details can be found in [3].
 - i. Cyclical
 - ii. Monotonic
 - iii. Without KL annealing strategy

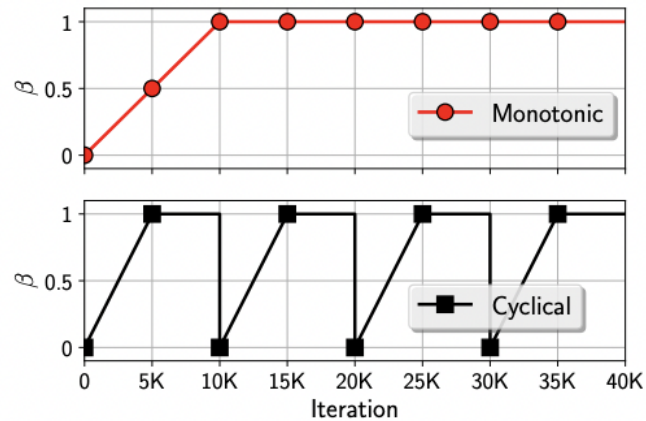


Figure 5. weight in different KL annealing strategies

5. Dataset

a. Training dataset

- i. train_img: 23410 png files
- ii. train_label: 23410 png files

b. Valadition dataset

- i. val_img: 630 png files
- ii. val_label: 630 png files

c. Testing dataset

- i. 5 video sequences are given. Each video sequence contains one first frame and 630 label frames.

6. Model Configurations

- a. Input images and labels will be resized to (32, 64) due to memory limitation.
- b. All modules e.g.
 - i. frame encoder
 - ii. pose encoder
 - iii. posterior generator
 - iv. decoder fusion
 - v. generator

are provided in directory named modules. **Please DO NOT change the structure of modules.**

c. Recommended command

- Training command

```
python Trainer.py --DR {YOUR_DATASET_PATH}  
--save_root {PATH_TO_SAVE_YOUR_CHECKPOINT}  
--fast_train
```

- --fast_train: is use fewer dataset and large learning rate to speed up your training

- Testing command

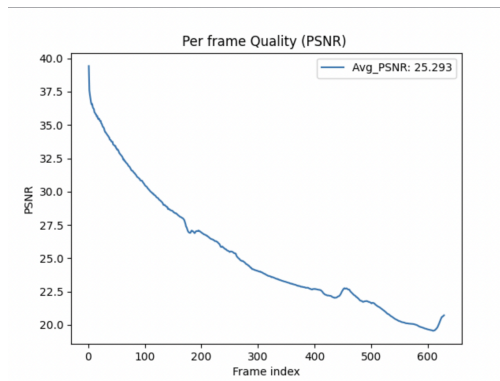
```
python Tester.py --DR {YOUR_DATASET_PATH}  
--save_root {PATH_TO_SAVE_YOUR_CHECKPOINT}  
--ckpt_path {PATH_TO_YOUR_CHECKPOINT}
```

7. Scoring Criteria

a. Report (60%)

- i. Derivate conditional VAE formula (5%)
- ii. Introduction (5%)
- iii. Implementation details (25%)
 1. How do you write your training/testing protocol (10%)
 2. How do you implement reparameterization tricks (5%)
 3. How do you set your teacher forcing strategy (5%)
 4. How do you set your kl annealing ratio (5%)
- iv. Analysis & Discussion (25%)
 1. Plot Teacher forcing ratio (5%)
 - a. Analysis & compare with the loss curve

2. Plot the loss curve while training with different settings. **Analyze the difference between them (10%)**
 - a. With KL annealing (Monotonic)
 - b. With KL annealing (Cyclical)
 - c. Without KL annealing
3. Plot the PSNR-per frame diagram in validation dataset (5%)



4. Other training strategy analysis (Bonus) (5%)




- b. Demo (50%)
 - i. Questions (20%)
 - ii. Kaggle competition (30%)
 1. Pass baseline: 20
 2. Top 30: 25
 3. Top 10: 30

Note: -5 points for wrong team name

Note: ii. Kaggle competition includes code explanation

8. Files provided

- a. Code
 - i. LAB4.zip is given on Kaggle. Structure would look like the picture below

- ▼  Lab4_template
 - ▼  modules
 - <> `__init__.py`
 - <> `layers.py`
 - <> `modules.py`
 - <> `Tester.py`
 - <> `Trainer.py`
 - <> `dataloader.py`
 -  `requirements.txt`

9. Hints

Friendly notification

- a. Do your work early.
- b. Read the reference papers
- c. Dataloader is provided

10. Upload file format

- a. Lab 4 - Conditional VAE (code)

File name: LAB4_{studentID}_{YOUR_NAME}.zip

example: LAB4_0812226_李阿金.zip

Note: -5 for wrong file name

- LAB4_0812226_李阿金.zip

- Your code

- b. Lab 4 - Conditional VAE (report)

File name: LAB4_{studentID}_{YOUR_NAME}.zip

example: LAB4_0812226_李阿金.zip

Note: -5 for wrong file name

11. References

- [1] C. Chan, et al., "Everybody Dance Now," ICCV, 2019
- [2] E. Denton, et al., "Stochastic Video Generation with a Learned Prior," ICML, 2018
- [3] H. Fu, et al., "Cyclical Annealing Schedule: A Simple Approach to Mitigating KL Vanishing," NAACL 2019