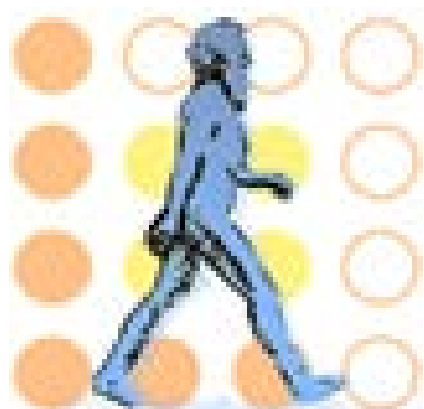




Deep Label Distribution Learning for Apparent Age Estimation

Xu Yang, Bin-Bin Gao, Chao Xing, Zeng-Wei Huo, Xiu-Shen Wei, Ying Zhou,
Jianxin Wu, and Xin Geng

ChaLearn Looking at People: Workshop and Competitions
@ICCV, 2015



presented by Bin-Bin Gao

Dec. 12, 2015 Santiago de Chile





Why is it difficult?

- It is difficult to provide **an exact answer**.

How old do these people look like?



A1: 30 **or** 32 years old;
A2: **Around** 31 years old;

.....

31 ± 4.24



A1: 18 **or** 20 years old;
A2: **May be** 20 years old;

.....

17 ± 1.93

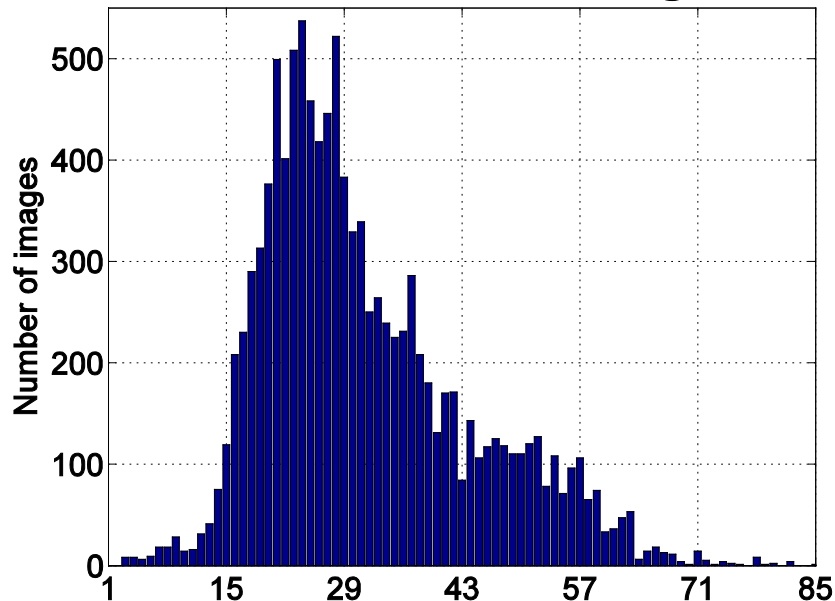


Why is it difficult?

- It is difficult to collect a **sufficient** and **complete** training dataset.

- ✓ **ChaLearn Competition**

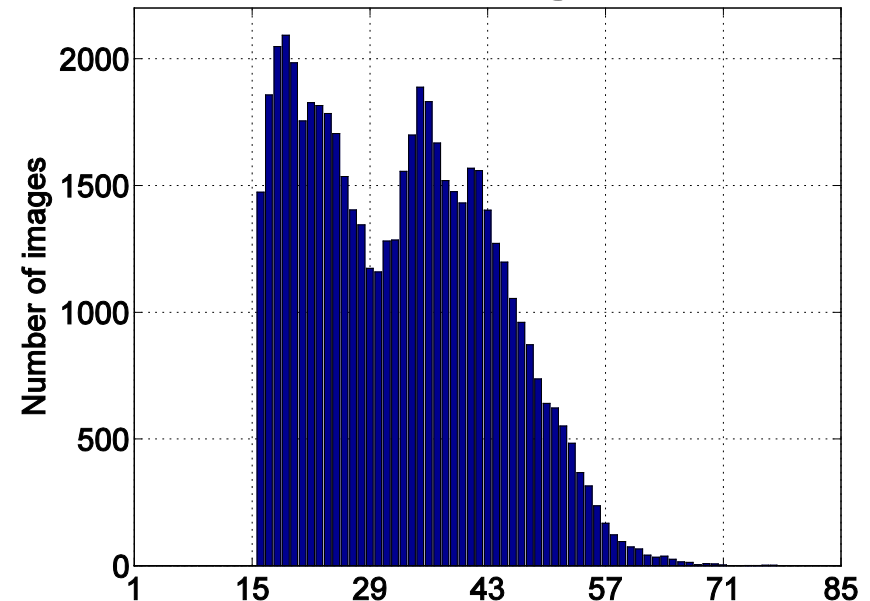
- Ages: from 1 to 85
- Training: 2,476 images
- Validation: 1,036 images



- ✓ Public age dataset

- Morph Album 2***

- Ages: from 15 to 77
- 55,134 images total



- Training data always has **small scale** and **imbalance**.

Some potential applications

- Although the task is very challenging, it has many **potential applications**.



Cigarette vending machine



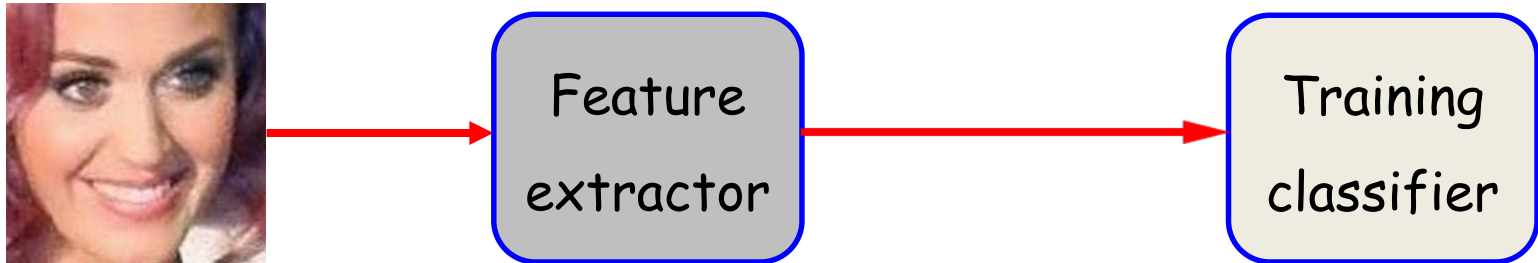
Kraft's vending machine

- Vending machines can prevent minors buying cigarettes and alcohol by estimating costumer's apparent age.



Many existing methods

✓ Hand-crafted feature



- BIF feature [Guo et al., CVPR 2009]
- OHRank [Chang et al., CVPR 2011]
- CCA, rCCA and kCCA [Guo et al. FGR 2013]
- IIS-LLD and CPNN [Gen et al., TPAMI 2013]
-



Many existing methods

✓ Deep learning



- Multi-scale CNN [Yi et al., ACCV 2014]
 - CNN based **regression** [Huerta et al., PRL 2015]
 - CNN for age group **classification** [Levi & Hassner, CVPR 2015]
 - DLA based on CNN features from different layers [Wang et al., WACV 2015]
 -
-



Motivations



30 ± 4.17



31 ± 4.24



32 ± 4.23



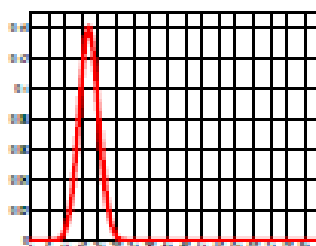
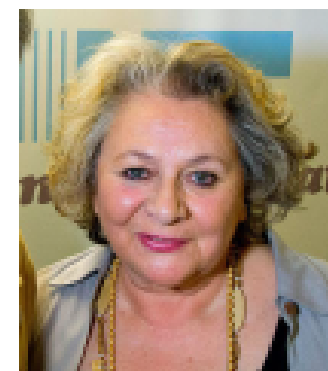
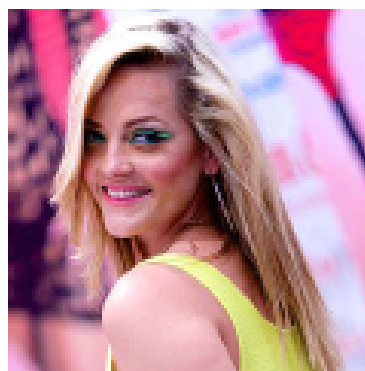
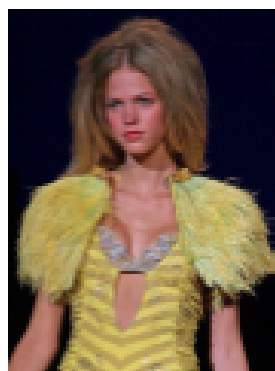
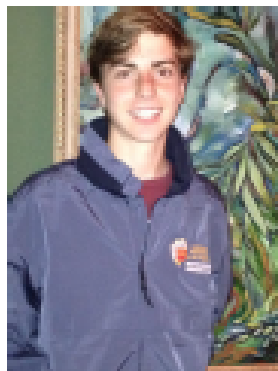
33 ± 2.01

Faces with similar ages look alike in terms of facial details such as wrinkles or skin smoothness. In other words, there is a *correlation among neighboring ages at both **image** and **feature** level.*

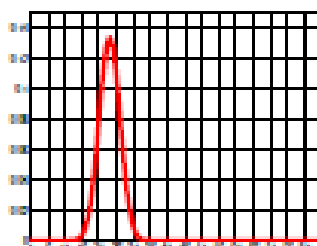
How to utilize the correlation?

✓ How to utilize the correlation?

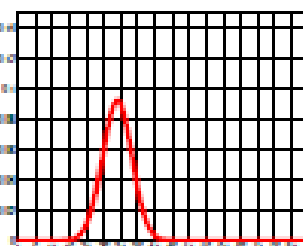
- Label Distribution (LD) Learning. But it does **not learn the visual representations**.



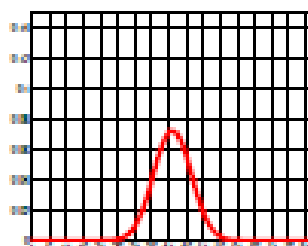
(a) (18,2.82)



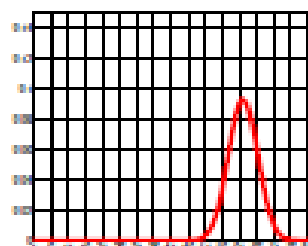
(b) (24,2.97)



(c) (30,4.30)



(d) (42,5.45)

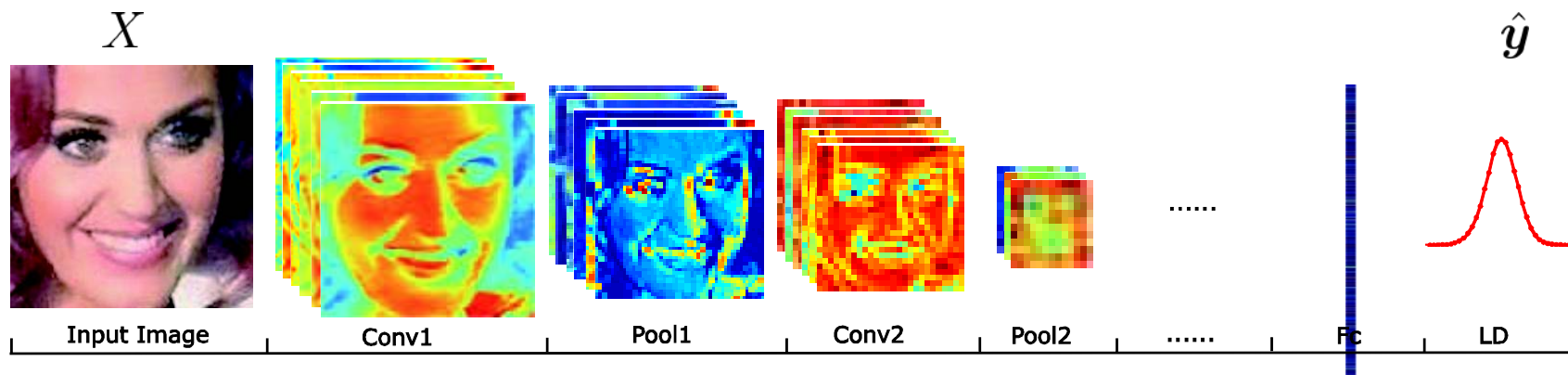


(e) (62,4.27)

- Generating LD $y_j = \frac{p(l_j|\mu, \sigma)}{\sum_k p(l_k|\mu, \sigma)}$, where $p(l_j|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(l_j - \mu)^2}{2\sigma^2}\right)$

Proposed methods

Deep Label Distribution Learning (DLDL)



✓ Formally:

- The goal of DLDL is to directly learn a conditional probability mass function $\hat{y} = p(y|X; \theta)$ from the training set, where θ is the parameter of the framework.

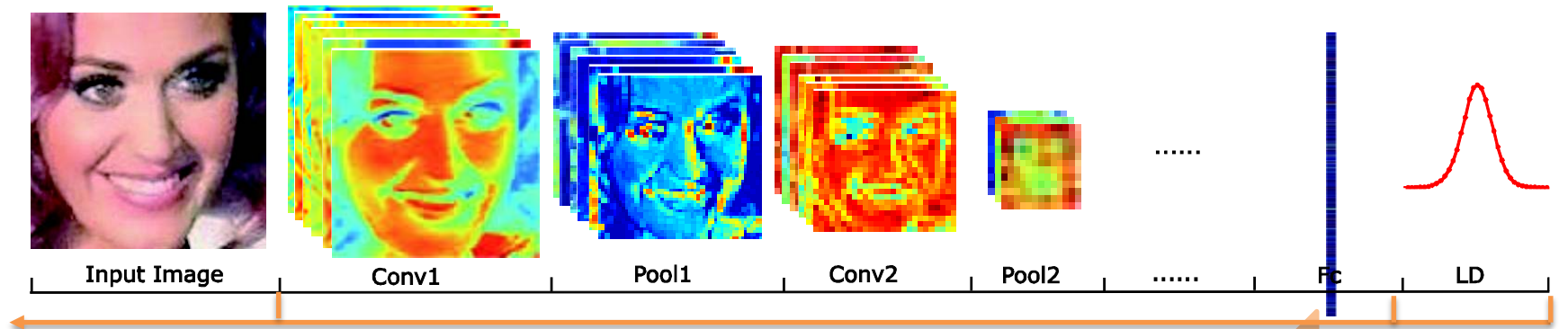
Proposed methods

Learning

✓ Forward:
 X

$$\hat{y}_j = \frac{\exp(x_j)}{\sum_t \exp(x_t)}$$

$$x = \phi(X; \theta) \quad \hat{y}$$



✓ Backward propagation:

- Objective function with K-L divergence:

$$\theta^* = \operatorname{argmin}_{\theta} \sum_k y_k \ln \frac{y_k}{\hat{y}_k} = \operatorname{argmin}_{\theta} \left[-\sum_k y_k \ln \hat{y}_k \right] T$$

- The gradient of the K-L loss function is given by

$$\frac{\partial T}{\partial x} = \hat{y} - y$$

$$\frac{\partial T}{\partial \theta} = \frac{\partial T}{\partial x} \times \frac{\partial x}{\partial \theta}$$

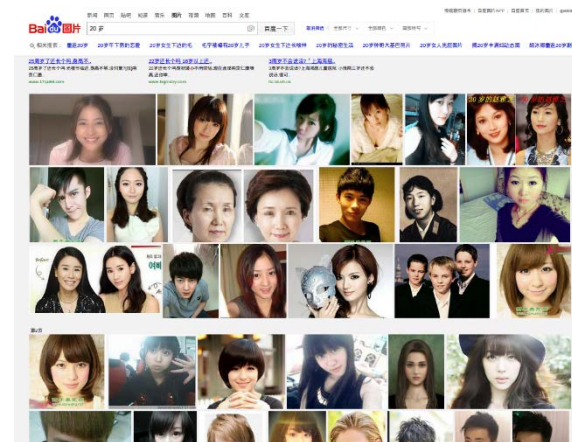
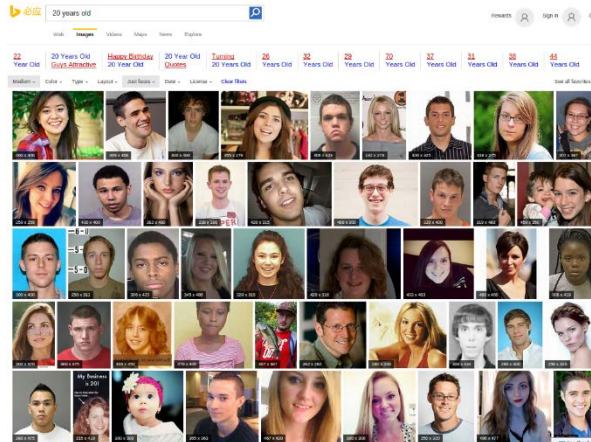
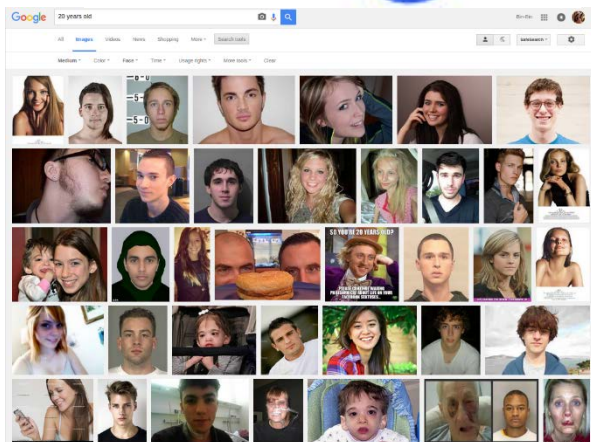
Our datasets

✓ Internet face images collecting

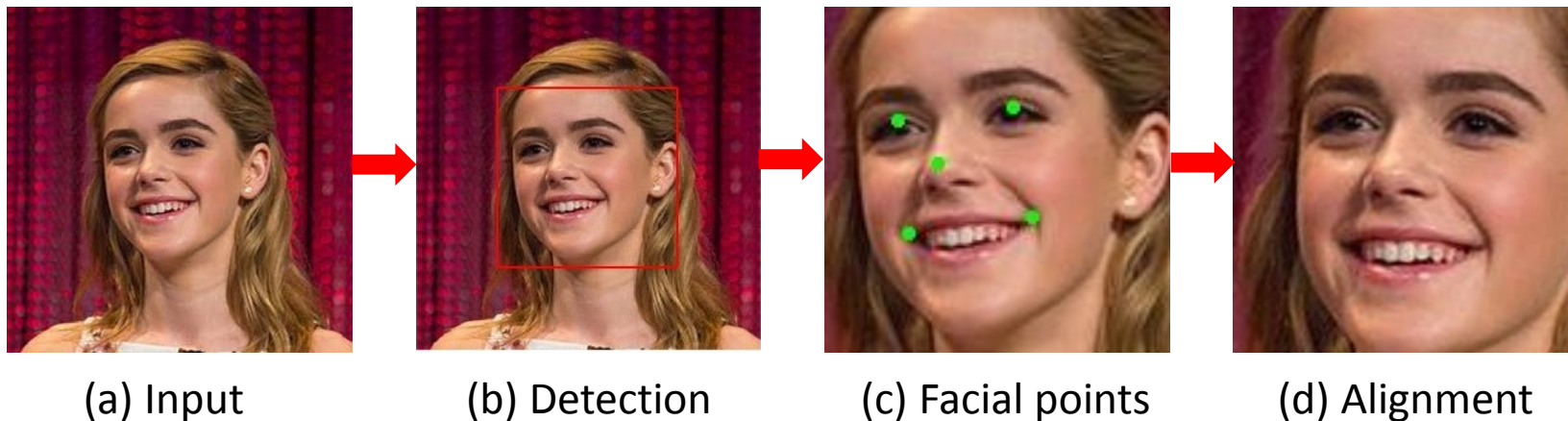
- A set of age related text enquires:
eg., "20 years old", "20th birthday" and "age-20" for the age of 20 years.
- We use **Google**, **Bing** and **Baidu** image search.

27,197 images

37,606 images



The face image pre-processing

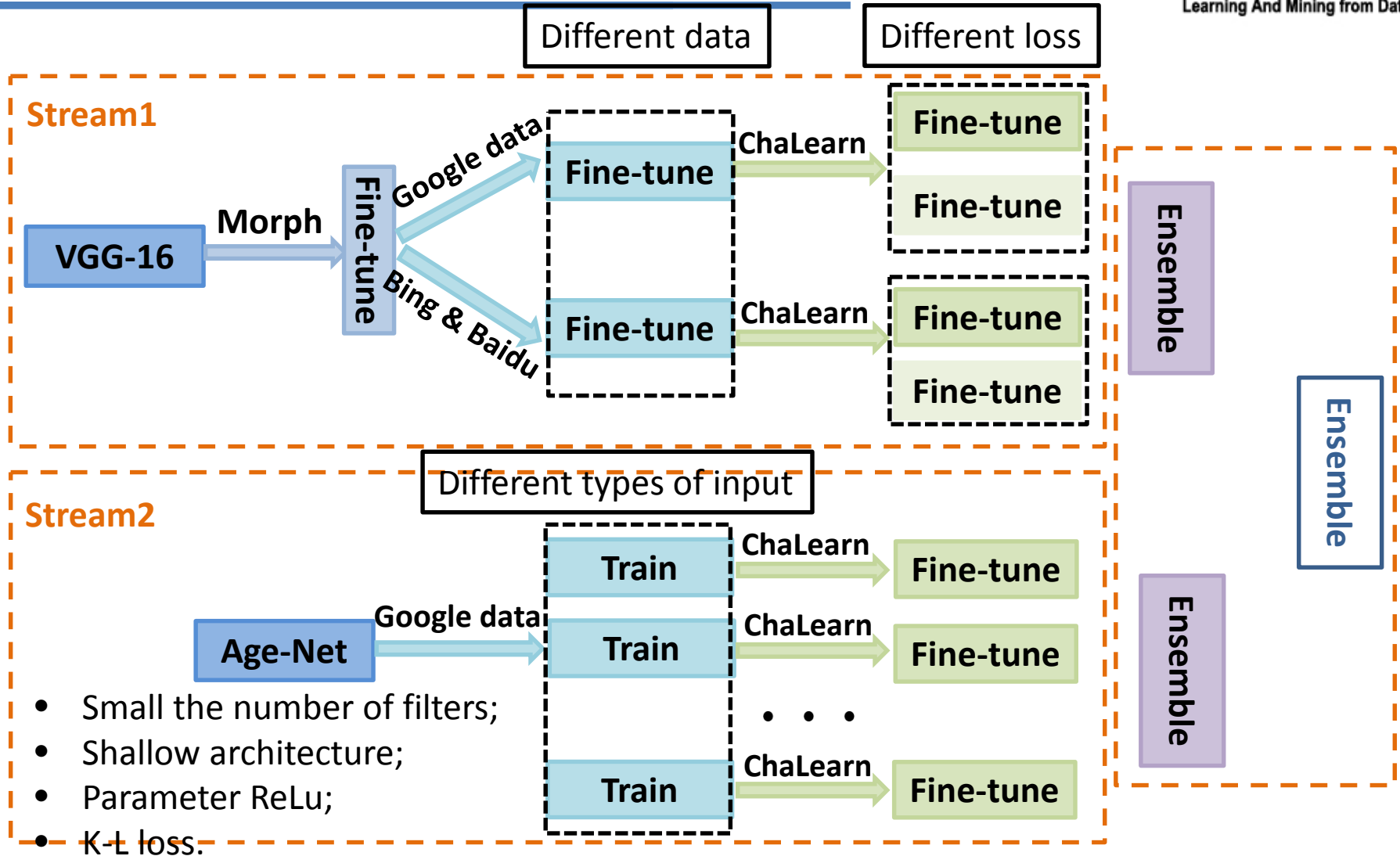


Three steps of the images pre-processing

- Face detection
- Facial points detection
- Face alignment



Model architecture





Training and prediction details

✓ Training

- Gaussian random initialization at different layers.

$1^{st} stream$: The last three layers \longrightarrow The last layer \longrightarrow The last layer.

$2^{st} stream$: All layers \longrightarrow The last layer.

✓ Prediction

- Different fusion strategy

Early fusion:

$1^{st} stream$: Prediction via measuring distance.

$2^{st} stream$: Averaging estimation distribution.

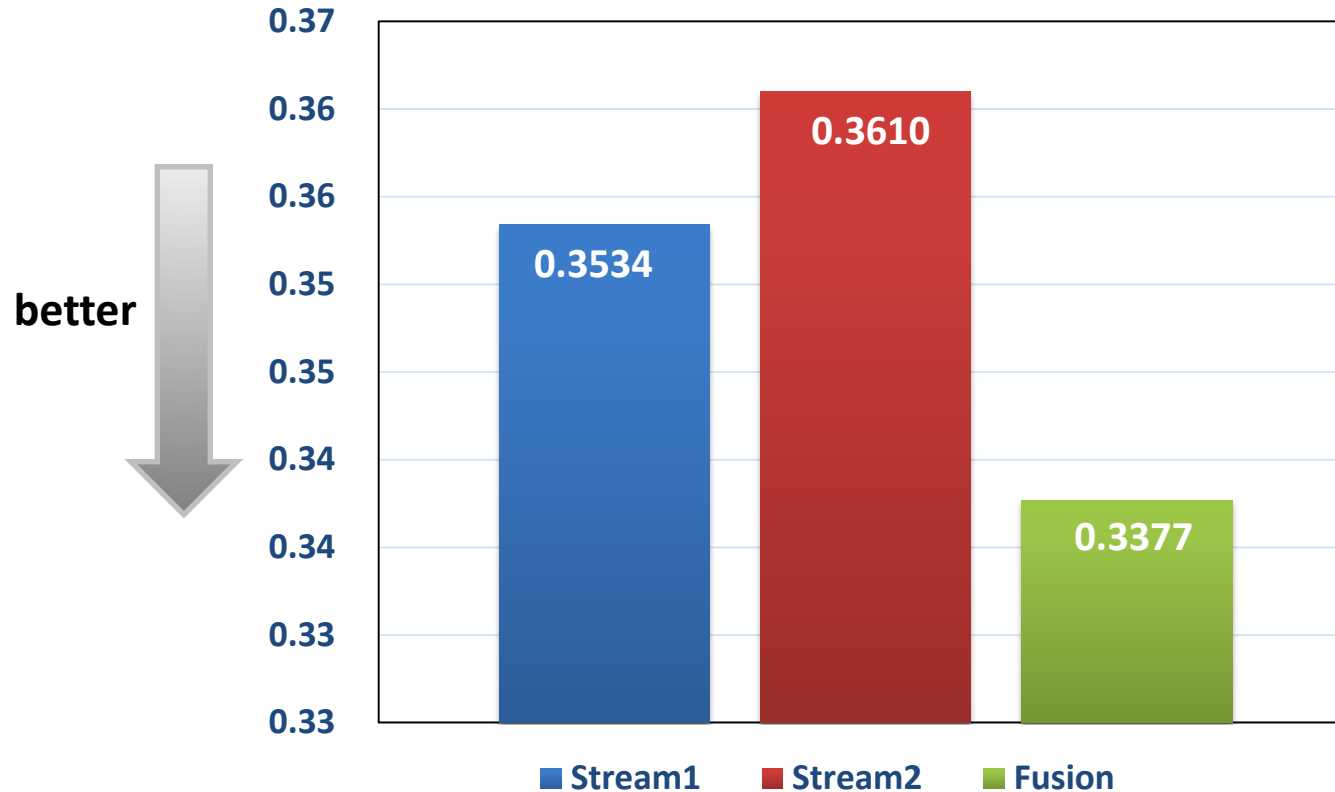
Late fusion :

Averaging the prediction age of two streams.



Comparison

Mean Error on Validation Set



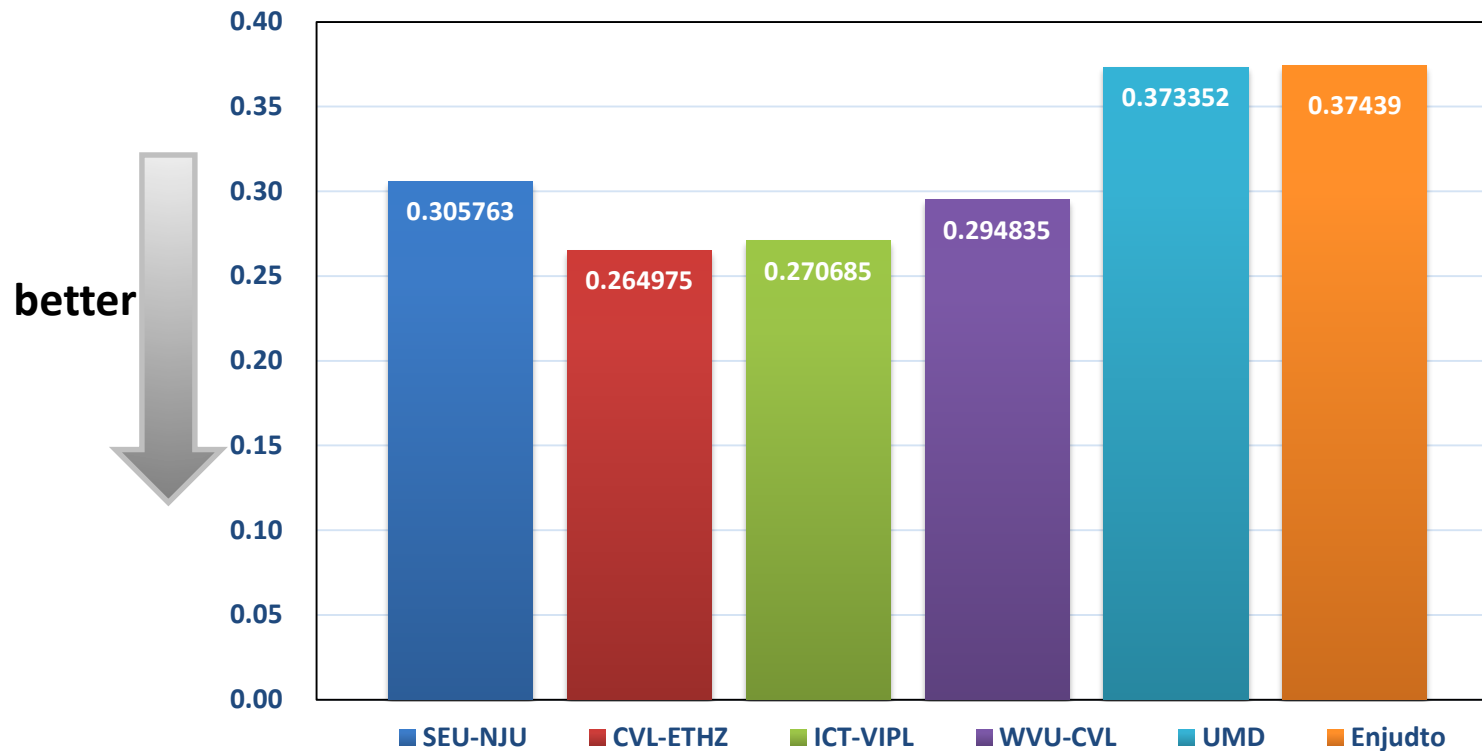
The fusion of the two stream is better than single stream.



Final results

The 4nd place with 0.3057 performance.

Mean Error on Test Set





Conclusions

- ✓ DLDL is an **end-to-end learning framework** which utilizes the **correlation** among neighboring labels in both feature learning and classifier learning;
 - ✓ DLDL can work when the training set is small.
 - ✓ **Ensemble strategy**: different dataset, different architecture, different initialization and different fusion.
-



Any questions

