

## Context / Problem

Facial expression recognition (FER) models lack interpretability

### FER classifier

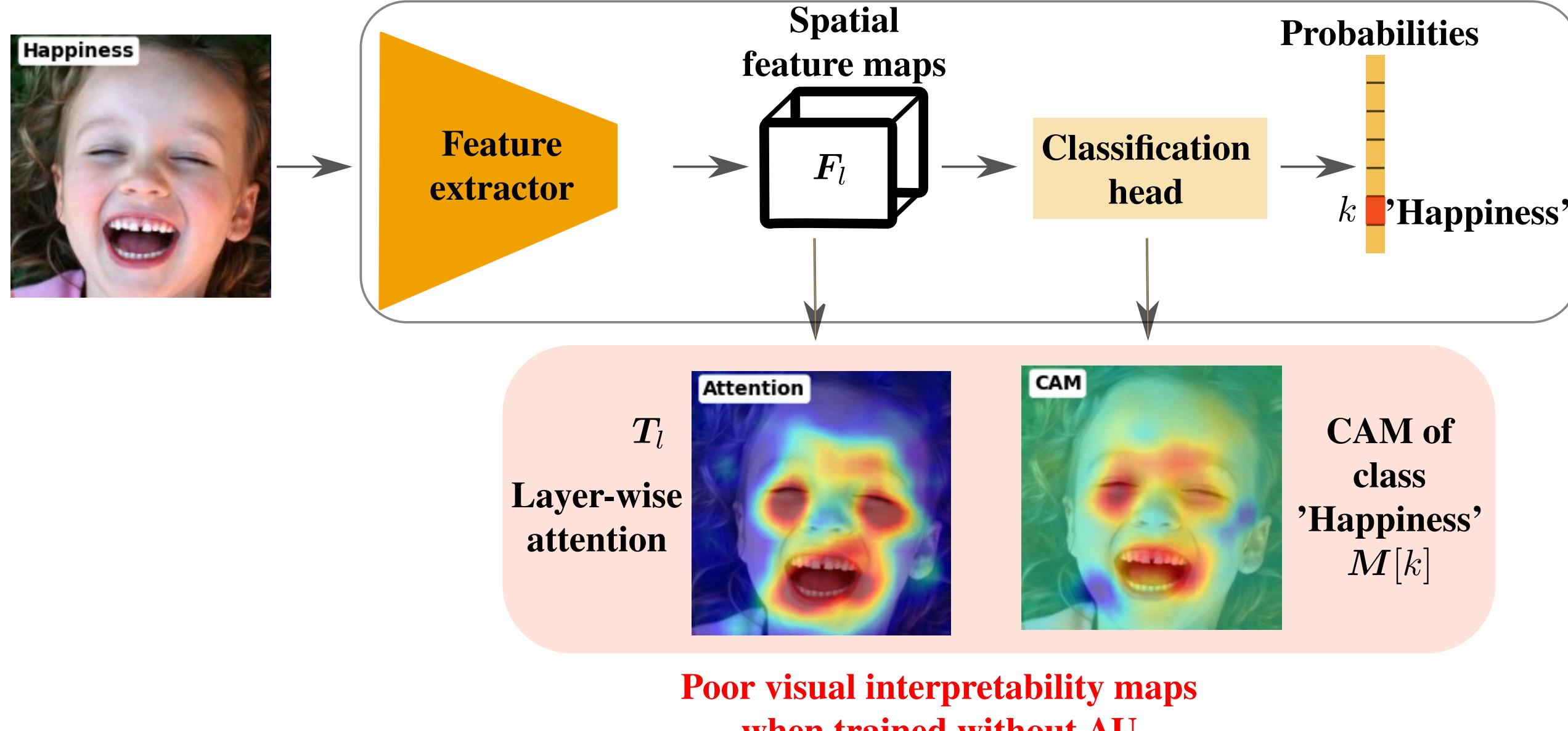


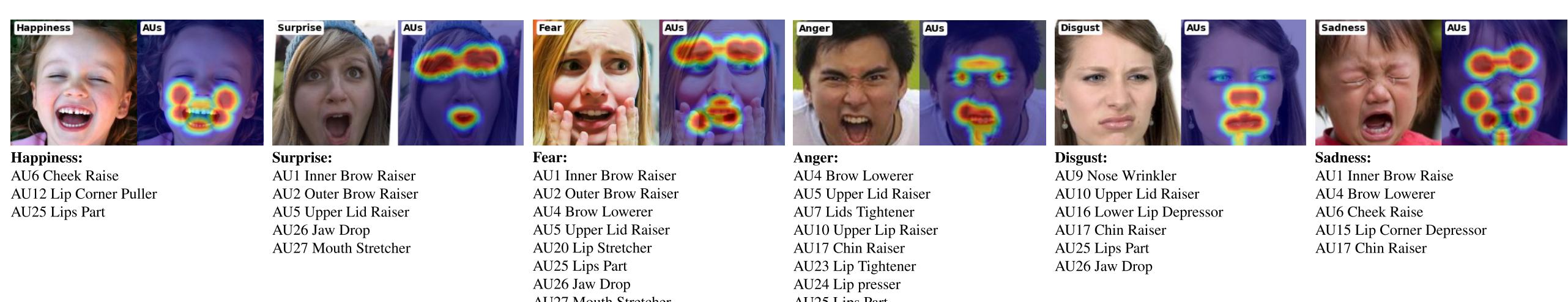
Figure 1. Poor interpretability of FER classifiers.

Our goal: Design an interpretable deep classifier for FER

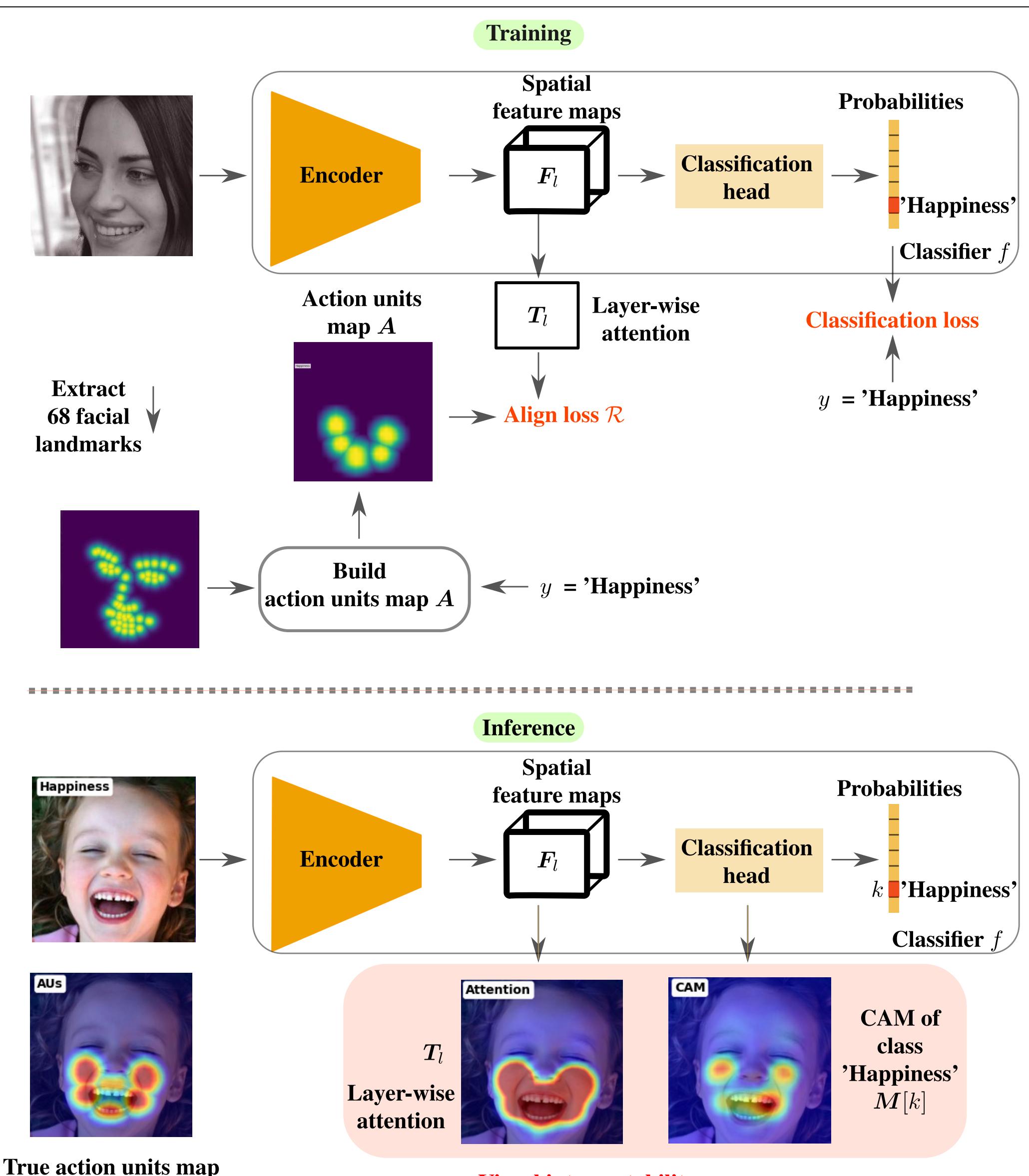
- Using similar cues used by experts: Facial Action Units (AUs)
- Allow the model to provide visual interpretability
- Maintain high classification performance

### Proposal:

- Build AU heatmaps for training: Requires only image class label
- Align layer-wise attention with AU heatmaps



### Proposal: Architecture



### Proposal: Training Loss

Total loss: classification, and spatial AU alignment,

$$\min_{\theta} -\log(f(\mathbf{X}; \theta)_y) + \lambda(1 - \mathcal{R}(\mathbf{T}_l, \mathbf{A})), \quad (1)$$

where,

$$\mathcal{R}(\mathbf{T}_l, \mathbf{A}) = \frac{\sum (\mathbf{T}_l \odot \mathbf{A})}{\|\mathbf{T}_l\|_2 \|\mathbf{A}\|_2}, \quad (2)$$

is the cosine similarity between the AU map  $\mathbf{A}$  and the layer-wise attention  $\mathbf{T}_l$  that is computed as,

$$\mathbf{T}_l = \frac{1}{n} \sum_{j=0}^n \mathbf{F}_l[j], \quad (3)$$

over the layer  $l$  and its spatial features  $\mathbf{F}_l$ .

## Empirical Results: RAF-DB & AffectNet Datasets

Classification (CL) and CAM-localization (CAM-COS) performance on RAF-DB and AffectNet test sets with and without AUs across methods:

Method	RAF-DB				AffectNet			
	CL		CAM-COS		CL		CAM-COS	
	w/o AU	w/ AU	w/o AU	w/ AU	w/o AU	w/ AU	w/o AU	w/ AU
CNN-based								
CAM (cvpr,2016)	88.20	<b>88.95</b>	0.55	<b>0.70</b>	60.88	<b>62.37</b>	0.56	<b>0.69</b>
WILDCAT (cvpr,2017)	88.26	<b>88.85</b>	0.52	<b>0.69</b>	59.88	<b>61.62</b>	0.62	<b>0.80</b>
GradCAM (iccv,2017)	88.39	<b>88.85</b>	0.55	<b>0.74</b>	60.77	<b>62.08</b>	0.53	<b>0.75</b>
GradCAM++ (wacv,2018)	87.84	<b>89.14</b>	0.60	<b>0.82</b>	60.22	<b>62.45</b>	0.66	<b>0.83</b>
ACoL (cvpr,2018)	87.94	<b>88.68</b>	0.54	<b>0.67</b>	58.28	<b>61.48</b>	0.55	<b>0.65</b>
PRM (cvpr,2018)	88.13	<b>88.88</b>	0.48	<b>0.59</b>	57.77	<b>60.97</b>	0.52	<b>0.75</b>
ADL (cvpr,2019)	87.45	<b>88.65</b>	0.50	<b>0.63</b>	57.88	<b>61.25</b>	0.54	<b>0.66</b>
CutMix (eccv,2019)	88.39	<b>88.59</b>	0.55	<b>0.57</b>	58.74	<b>59.88</b>	0.56	<b>0.58</b>
LayerCAM (ieee,2021)	87.90	<b>88.88</b>	0.60	<b>0.84</b>	60.77	<b>62.45</b>	0.66	<b>0.83</b>
Transformer-based								
TS-CAM (iccv,2021)	86.70	<b>88.00</b>	0.58	<b>0.71</b>	58.99	<b>59.54</b>	0.57	<b>0.58</b>
APViT (ieee,2022)	91.00	<b>91.03</b>	--	--	60.62	<b>62.28</b>	--	--

Attention-localization (ATT-COS) (at layer 5) performance over RAF-DB and AffectNet test sets with and without AUs:

Methods / Case	RAF-DB		AffectNet	
	w/o AU	w/ AU	w/o AU	w/ AU
<b>CNN-based</b>				
CAM(cvpr,2016)	0.57	<b>0.85</b>	0.64	<b>0.82</b>
WILDCAT (cvpr,2017)	0.47	<b>0.85</b>	0.61	<b>0.81</b>
GradCAM (iccv,2017)	0.63	<b>0.85</b>	0.65	<b>0.82</b>
GradCAM++ (wacv,2018)	0.52	<b>0.87</b>	0.65	<b>0.82</b>
ACoL (cvpr,2018)	0.46	<b>0.84</b>	0.60	<b>0.81</b>
PRM (cvpr,2018)	0.43	<b>0.85</b>	0.55	<b>0.82</b>
ADL (cvpr,2019)	0.51	<b>0.85</b>	0.65	<b>0.83</b>
CutMix (eccv,2019)	0.51	<b>0.80</b>	0.57	<b>0.82</b>
LayerCAM (ieee,2021)	0.52	<b>0.86</b>	0.65	<b>0.82</b>
<b>Transformer-based</b>				
TS-CAM (iccv,2021)	0.55	<b>0.88</b>	0.48	<b>0.79</b>
APViT (ieee,2022)	0.38	<b>0.85</b>	0.45	<b>0.84</b>

Visual Results: Better interpretability

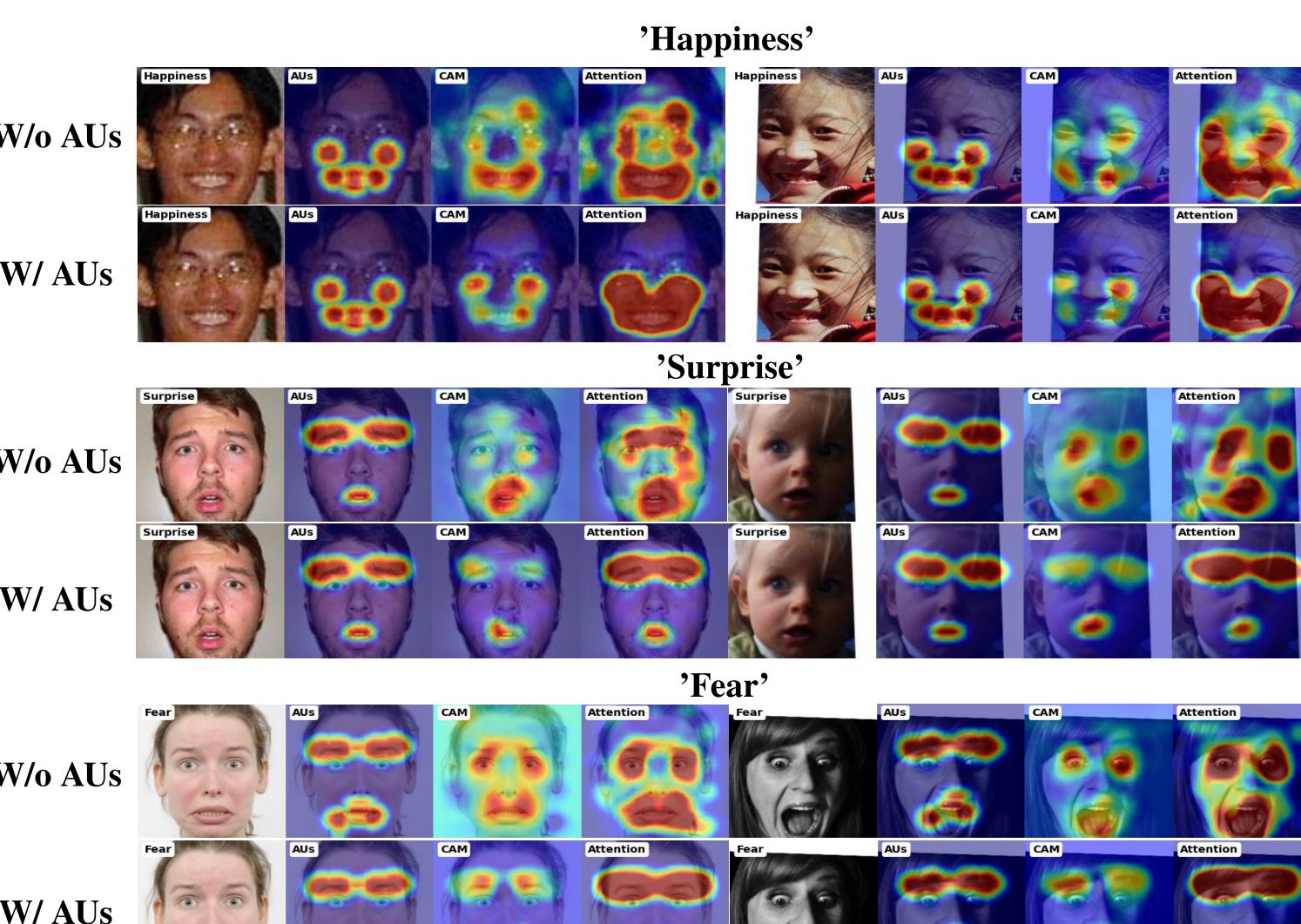


Figure 2. Illustration of interpretability prediction over RAF-DB test samples using CAM method with and without action units alignment. From left to right: Input image, true action units map  $\mathbf{A}$ , CAM  $M[k]$ , attention  $T_5$ , 'Happiness', 'Surprise', 'Fear', 'Anger', 'Disgust', 'Sadness'.

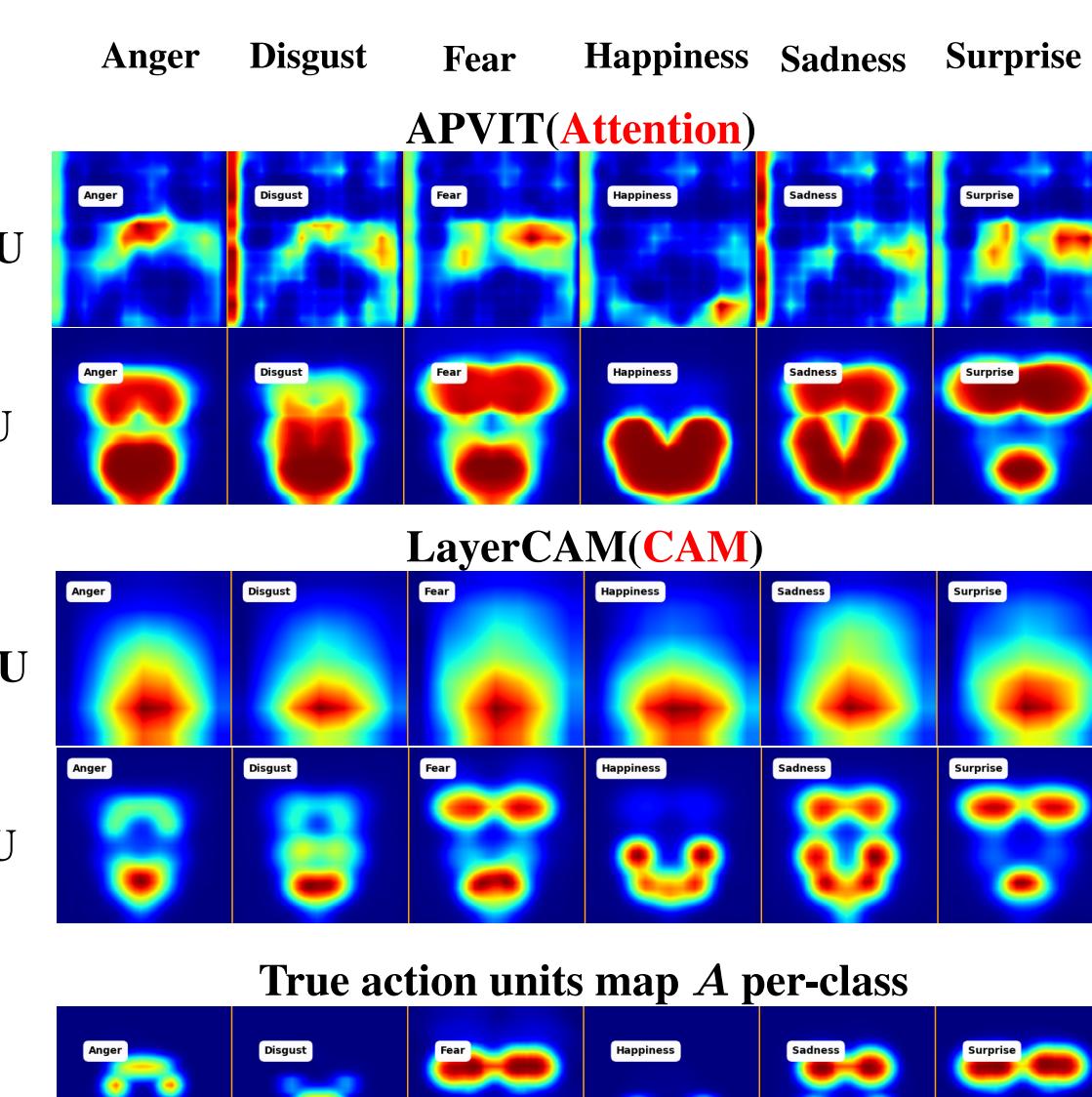


Figure 3. Illustration of per-class average attention and CAM maps over all test set of RAF-DB with and without AU alignment. Expressions from left to right: 'Anger', 'Disgust', 'Fear', 'Happiness', 'Sadness', 'Surprise'.



Figure 4. Ablation on the RAF-DB test set: impact of  $\lambda$  over classification and localization (interpretability) performance.