# Learning Facial Action Units from Web Images with Scalable Weakly-supervised Spectral Clustering <sup>1</sup>Beijing University of Posts and Telecomm. <sup>2</sup>Carnegie Mellon University <sup>3</sup>The Ohio State University



## Problem

## ▲ Facial Action Unit (AU) detection



AU6 Cheek raiser AU12

Inner brow raisor

Lip corner depressor

Happiness

## Motivation

- 1. Utilize large and freely available web images
- 2. Avoid manual annotation laborious and error-prone
- **3.** Improve model performance with free unannotated data



## Weakly-supervised Spectral Clustering (WSC)

Step 1. Weakly supervised embedding (WSE) Step 2. Re-annotation via rank-order clustering



## ▲ Alternative methods

Methods	UD	PN	SL	IE
STM, CPM [1]		$\mathbf{\otimes}$	$\mathbf{x}$	
GFK, LapSVM [2]		$\mathbf{\mathbf{S}}$	$\mathbf{x}$	
Spectral/K-means clu.		$\mathbf{x}$	$\mathbf{x}$	
WSC				

**UD**: Unannotated data, **PN**: Pruning noises, SL: Scalability, IE: Identity exemption

- [1] "Selective transfer machine for personalized facial action unit detection," in CVPR, 2013.
- [2] "Geodesic flow kernel for unsupervised domain adaptation," in CVPR, 2012.



Step 1: Weakly supervised embedding

**Re-annotation:** pos, neg



Step 2: Re-annotation via clustering

## <sup>1</sup>Kaili Zhao, <sup>2</sup>Wen-Sheng Chu, <sup>3</sup>Aleix M. Martinez

visual similarity and weak annotation in 1 million images

$$\min_{\mathbf{W}\in\mathbb{R}^{N\times K}} \underbrace{f(\mathbf{W},\mathbf{L})}_{\mathbf{W}\in\mathbb{R}^{N\times K}} + \frac{\lambda}{|\mathcal{G}|} \underbrace{\psi(\mathbf{W},\mathcal{G})}_{\mathbf{W}\in\mathbb{R}^{N\times K}} \quad \text{s.t.} \quad \mathbf{W}^{\top}\mathbf{W} = \mathbf{I}_{K}$$

• Visual similarity:

$$f(\mathbf{W}, \mathbf{L}) = \operatorname{Tr}(\mathbf{W}^{\top} \mathbf{L} \mathbf{W}), \ \mathbf{L} = \mathbf{D} - \mathbf{A}, \ A_{ij} = \begin{cases} \exp(-\gamma d(\mathbf{x}_i, \mathbf{x}_j)), \\ 0 \end{cases}$$

• Weak annotation: (Agreement)

$$(\mathbf{W}, \mathcal{G}) = \frac{1}{n_g} \sum_{\mathbf{w}_i \in \mathcal{G}_g} (\mathbf{w}_i - \overline{\mathbf{w}}_g)^\top (\mathbf{w}_i - \overline{\mathbf{w}}_g) = \frac{1}{n_g} \sum_{\mathbf{w}_i \in \mathcal{G}_g} \operatorname{Tr}(\mathbf{W}_i)$$

expansion and group decomposition

• Analytical solution: 
$$\mathbf{W}_g^{\star} = (\mathbf{I}_{n_g} + \frac{2\tilde{\lambda}}{n_g}\mathbf{C}_g)^{-1}\mathbf{V}_g \leftarrow \text{Inverse is slow and numerical solution}$$

• 10x-Faster solution: 
$$\mathbf{W}_i = \frac{1}{a} \mathbf{V}_i - \frac{b}{a(a+bn_g)} \sum_j \mathbf{V}_j, \ a = 1 + \frac{2\lambda}{n_g}, b = \frac{2\lambda}{-n_g^2}$$

### Algorithm 1 Weakly Supervised Spectral Embedding Algorithm 2 Stochastic Spectral Embedding **Input:** Laplacian matrix $\mathbf{L} \in \mathbb{R}^{N \times N}$ , orthonormal matrix $\mathbf{W}_0$ $\mathbb{R}^{N \times K}$ , stepsize $\eta$ , update ratio $\gamma$ , and tuning parameter $\lambda$ **Output:** An orthonormal matrix $\mathbf{W} \in \mathbb{R}^{N \times K}$ $\eta$ , update ratio $\gamma$ , and tuning parameter $\lambda$ 1: $a_0 = 1, t = 0$ **Output:** An orthonormal matrix $\mathbf{W} \in \mathbb{R}^{N \times K}$ : while not converge do while $t \leq T$ do if $f(\mathbf{W}_t) + \lambda \psi(\mathbf{W}_t, \mathcal{G}) \ge Q_L(\mathbf{W}_t, \mathbf{V})$ then for b = 1, ..., B do $\eta = \gamma \eta$ end if $\mathbf{V} = \mathbf{W}_t - \eta(2\mathbf{L}\mathbf{W}_t)$ for $\mathcal{G}_g \in \mathcal{G}$ do $\mathbf{W}_g = (\mathbf{I}_{n_g} + \frac{2\lambda}{n_g} \mathbf{C}_g)^{-1} \mathbf{V}_g$ // Update each group of $\mathbf{W}$ end for end while end for 8: $\mathbf{W} = \mathbf{W}_t$ learning problem," in ICDM, 2009. 13: end while 14: $W = W_t$ of images," in ECCV, 2014.







## Experiments

## ▲ EmotioNet dataset [6]

- 1M web images: 50K images (5%) were manually labeled by experts.
- 7 AUs with base rate > 5% were chosen for experiments.

### ▲ Settings

- 25K/25K partition of labeled images for training/test (following [6])
- Weak annotations were obtained by an AlexNet pre-trained on BP4D.
- The remaining 950K demonstrates the use of unlabeled images.

### ▲ Comparisons

• Annotations: **wlb:** weak annotation wsc: our annotation

gt: human annotation

• Experiment 1: {gt25k} + {[**wlb** | **wsc** | **gt**]10k}

### Findings:

- wsc >> wlb (> 20%)
- wsc ~= gt (~4%)
- Experiment (2): AlexNet vs DRML [7] Findings:
- DRML >~ AlexNet
- More unlabeled data plus WSC helps both
- Experiment 3: SSL methods

### Findings:

- SSL suffers from noisy data due to smoothness
- LapSVM is slow and fails to scale up
- WSC is efficient and best performer
- Experiment 4: Large-scale evaluation Findings:
- WSC scales to 1M!
- More improvement with more WSCannotated images

AU	AlexNet $\begin{cases} gt15k \\ wlb10k \end{cases}$	$\begin{cases} AlexNe \\ gt15k \\ wsc10k \end{cases}$	$ \left\{ \begin{array}{c} \text{dexNet} \\ \text{gt25k} \end{array} \right\} $	${ {gt25k } \\ }$	$\begin{cases} AlexNet \\ gt25k \\ wsc25k \end{cases}$	$DRML \\ \begin{cases} gt25k \\ wsc25k \end{cases}$
1	11.8	19.8	24.2	25.3	25.3	[26.3]
4	23.9	32.5	34.7	[35.7]	34.5	35.5
5	26.6	37.6	39.5	40.0	39.3	[40.3]
6	58.8	73.5	73.1	75.3	75.6	[78.7]
12	82.1	87.1	86.8	86.6	87.4	[88.1]
25	82.1	84.3	88.5	[88.9]	88.8	[88.9]
26	24.3	40.2	45.6	46.2	47.7	[49.1]
Avg.	44.2	53.6	56.1	56.9	57.0	[58.1]

3	F1					S score				
ATI	GF ∫gt25	K Laj 5k}∫∫§	pSVM gt25k ∖	TSVM ∫ gt25k	[ ]	G ∫gt	FK 25k∖	Laj ∫	pSVM gt25k ∖	TSVM ∫ gt25k ∖
	)	∫	$1b10k\int$	ulb25k	<u>}</u>	l	5	lι	$1b10k \int$	$\left\{ \text{ulb25k} \right\}$
1	19.3	1.	2	24.1		66	.1	82	2.3	70.2
4	31.0	31.0 25.7		32.3		61.1		85.3		62.5
5	31.8	23	.1	40.3		61	.1	60	).7	80.6
6	73.8	58	.3	75.7		71	.7	70	0.0	79.1
12	85.1	57	.7	87.4		75	5.5 50.9		).9	80.2
25	85.8	88	.9	88.2		72	2.4 79.4		9.4	78.5
26	39.0	5.	0	47.0		69	69.5 83.2		3.2	78.4
Avg.	52.2	37	.0	56.4		68.2 73.1		3.1	75.6	
4	2	20k   200k			400k			1M		
AU	wlb	WSC	wlb	wsc		wlb	WSC		wlb	wsc
1	17.6	18.3	17.8	19.3		16.9	[21.3	5]	17.6	21.2
4	20.3	20.5	19.0	20.4		18.9	21.3	6	18.4	[22.1]
5	28.5	28.9	30.1	30.8		31.5	33.4		30.8	[41.6]
6	72.4	74.1	75.9	76.9	,	76.3	78.6		77.4	[79.3]
12	76.7	85.8	79.1	86.4	,	79.3	87.8		81.4	[88.2]
25	84.7	85.7	85.4	85.9		79.4	86.1	-	86.1	[89.1]
26	32.4	34.9	32.7	36.0		33.3	36.1		33.3	[47.2]

- Avg. 47.5 49.7 | 48.5 50.8 | 47.9 52.1 | 49.3 [55.5] [6] "EmotioNet challenge: Recognition of facial expressions of emotion in the wild," in CVPRW, 2017.
- 7] "Deep region and multi-label learning for facial action unit detection," in CVPR, 2016.