

Learning Facial Action Units from Web Images with Scalable Weakly-supervised Spectral Clustering

¹Beijing University of Posts and Telecomm. ²Carnegie Mellon University ³The Ohio State University¹Kaili Zhao, ²Wen-Sheng Chu, ³Aleix M. Martinez

Problem

▲ Facial Action Unit (AU) detection

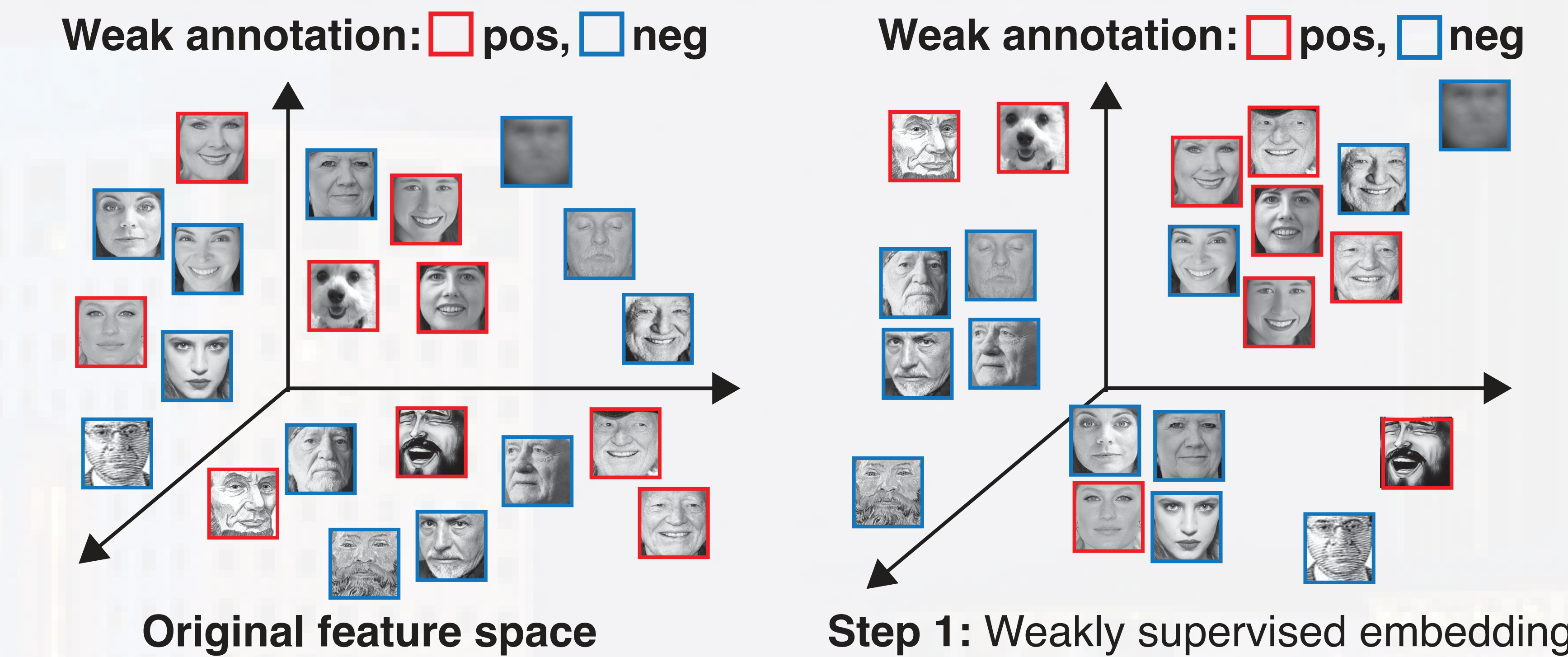


▲ Motivation

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]	✓	✗	✗	✗
GFK, LapSVM [2]	✓	✗	✗	✓
Spectral/K-means clu.	✓	✗	✗	✓
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.

Scalable weakly-supervised spectral embedding

▲ **Objective:** Learn an embedding space with coherence among **visual similarity** and **weak annotation** in 1 million images

- **Formulation:**
$$\min_{\mathbf{W} \in \mathbb{R}^{N \times K}} \underbrace{f(\mathbf{W}, \mathbf{L})}_{\text{Visual similarity}} + \frac{\lambda}{|\mathcal{G}|} \underbrace{\psi(\mathbf{W}, \mathcal{G})}_{\text{Weak annotation}} \quad \text{s.t.} \quad \mathbf{W}^\top \mathbf{W} = \mathbf{I}_K$$
- **Visual similarity:** $f(\mathbf{W}, \mathbf{L}) = \text{Tr}(\mathbf{W}^\top \mathbf{L} \mathbf{W}), \mathbf{L} = \mathbf{D} - \mathbf{A}, A_{ij} = \begin{cases} \exp(-\gamma d(\mathbf{x}_i, \mathbf{x}_j)), & \text{if } \mathbf{x}_i \in \mathcal{N}_k(\mathbf{x}_j) \\ 0 & \text{otherwise} \end{cases}$ (Spectral clustering)
- **Weak annotation:** $\psi(\mathbf{W}, \mathcal{G}) = \frac{1}{n_g} \sum_{\mathbf{w}_i \in \mathcal{G}_g} (\mathbf{w}_i - \bar{\mathbf{w}}_g)^\top (\mathbf{w}_i - \bar{\mathbf{w}}_g) = \frac{1}{n_g} \sum_{\mathbf{w}_i \in \mathcal{G}_g} \text{Tr}(\mathbf{W}^\top \mathbf{C}_g \mathbf{W})$ (Agreement)

▲ **Solution:** Address the nonsmooth nature of $\psi(\mathbf{W}, \mathcal{G})$ using first-order Taylor expansion and group decomposition

- **Analytical solution:** $\mathbf{W}_g^* = (\mathbf{I}_{n_g} + \frac{2\tilde{\lambda}}{n_g} \mathbf{C}_g)^{-1} \mathbf{V}_g \leftarrow$ Inverse is slow and numerically unstable!
- **10x-Faster solution:** $\mathbf{W}_i = \frac{1}{a} \mathbf{V}_i - \frac{b}{a(a + bn_g)} \sum_j \mathbf{V}_j, a = 1 + \frac{2\tilde{\lambda}}{n_g}, b = \frac{2\tilde{\lambda}}{-n_g^2}$

▲ Optimization: Accelerated gradient descent + stochastic extension

Algorithm 1 Weakly Supervised Spectral Embedding

Input: Laplacian matrix $\mathbf{L} \in \mathbb{R}^{N \times N}$, orthonormal matrix $\mathbf{W}_0 \in \mathbb{R}^{N \times K}$, stepsize η , update ratio γ , and tuning parameter λ

Output: An orthonormal matrix $\mathbf{W} \in \mathbb{R}^{N \times K}$

```

1:  $a_0 = 1, t = 0$ 
2: while not converge do
3:   if  $f(\mathbf{W}_t) + \lambda \psi(\mathbf{W}_t, \mathcal{G}) \geq Q_L(\mathbf{W}_t, \mathbf{V})$  then
4:      $\eta = \gamma \eta$ 
5:   end if
6:    $\mathbf{V} = \mathbf{W}_t - \eta(2\mathbf{L}\mathbf{W}_t)$ 
7:   for  $\mathcal{G}_g \in \mathcal{G}$  do
8:      $\mathbf{W}_g = (\mathbf{I}_{n_g} + \frac{2\lambda}{n_g} \mathbf{C}_g)^{-1} \mathbf{V}_g$  // Update each group of  $\mathbf{W}$ 
9:   end for
10:   $a_t = \frac{2}{t+3}$ 
11:   $\mathbf{W}_t = \mathbf{W}_t + \frac{1-a_{t-1}}{a_{t-1}} \cdot a_t (\mathbf{W}_t - \mathbf{W}_{t-1})$ 
12:   $\mathbf{W}_t = \text{orth}(\mathbf{W}_t)$  // Enforce  $\mathbf{W}_t$  to be orthonormal
13: end while
14:  $\mathbf{W} = \mathbf{W}_t$ 

```

Algorithm 2 Stochastic Spectral Embedding

Input: Laplacian matrix $\mathbf{L} \in \mathbb{R}^{N \times N}$, orthonormal matrix $\mathbf{W}_0 \in \mathbb{R}^{N \times K}$, number of batches B , number of iterations T , stepsize η , update ratio γ , and tuning parameter λ

Output: An orthonormal matrix $\mathbf{W} \in \mathbb{R}^{N \times K}$

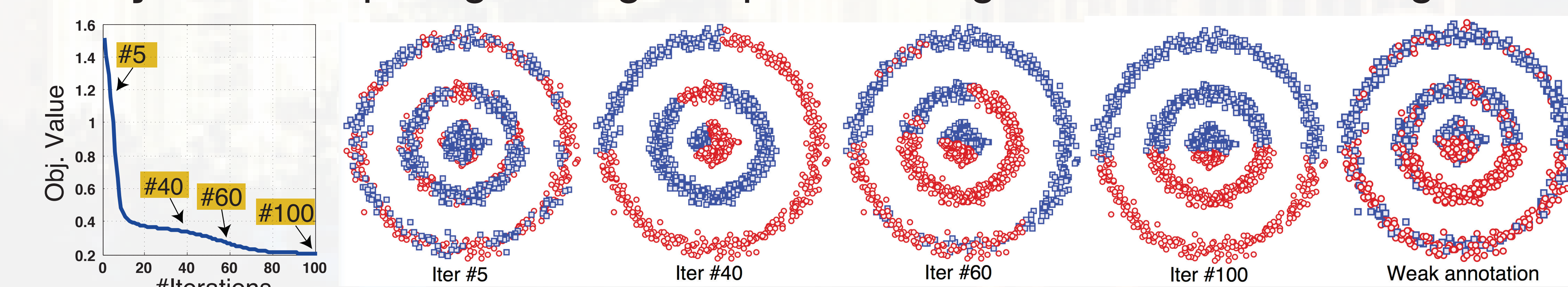
```

1: while  $t \leq T$  do
2:   for  $b = 1, \dots, B$  do
3:      $\tilde{\mathbf{L}}_t = \text{sampling}(\mathbf{L})$  // Perform edge sampling
4:     Solve  $\mathbf{W}_t$  using Algorithm 1 with  $(\tilde{\mathbf{L}}_t, \mathbf{W}_{t-1}, \eta, \gamma, \lambda)$ 
5:      $\mathbf{W}_t = \text{orth}(\mathbf{W}_t)$  // Enforce  $\mathbf{W}_t$  to be orthonormal
6:   end for
7: end while
8:  $\mathbf{W} = \mathbf{W}_t$ 

```

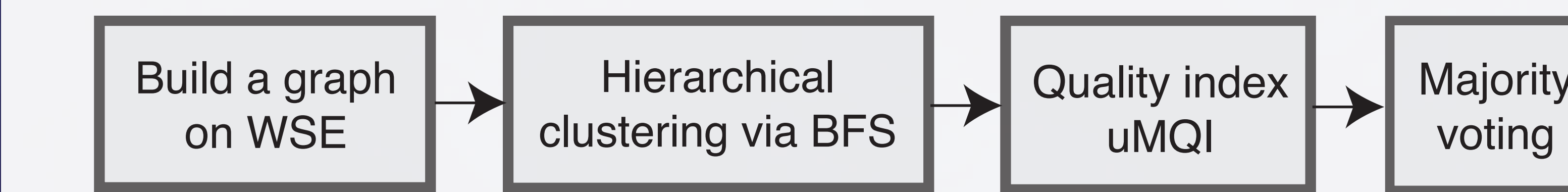
- [3] "Accelerated gradient method for multi-task sparse learning problem," in ICDM, 2009.
 [4] "Spectral clustering with a convex regularizer on millions of images," in ECCV, 2014.

▲ Toy EX: Group neighboring samples with high weak annotation agreement



Re-annotation via clustering

▲ Re-annotation pipeline



▲ Build a graph based on WSE embedding

- **Rank-order distance** [5]: Same-class samples have similar neighbors
- Standard L1/L2 distance is sensitive to biased distribution of AU data

▲ Hierarchical clustering

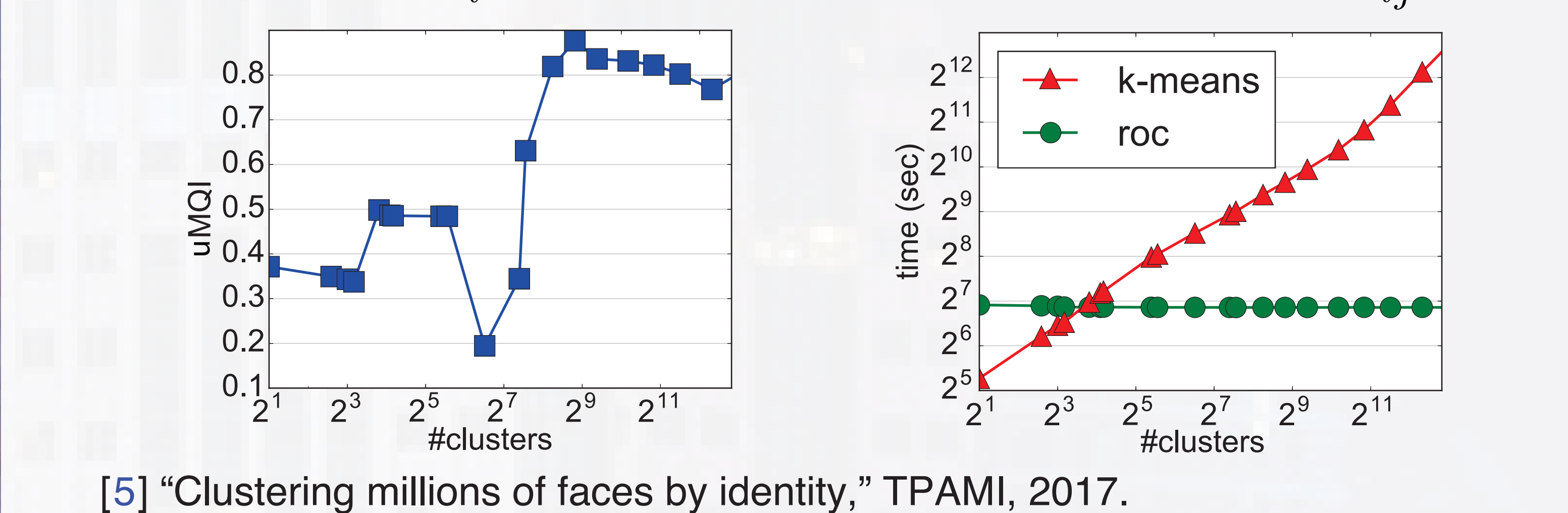
- **Scalable to 1M images:** $\mathcal{O}(Nk) + \mathcal{O}(N)$ vs k-means $\mathcal{O}(tKNd)$
- **Noise/outlier pruning** by identifying clusters with rare samples
- Avoid the non-trivial **#cluster** as input

▲ Undirected modularization quality index (uMQI)

• **Intra-cluster connectivity** vs **inter-cluster isolation**

$$\text{uMQI}(\mathcal{C}) = \underbrace{\mathbb{E}[\text{intra}(\mathcal{C})]}_{\text{Intra-cluster connectivity}} - \underbrace{\mathbb{E}[\text{inter}(\mathcal{C})]}_{\text{Inter-cluster isolation}}, \mathcal{C} = \{\mathcal{C}_i, \forall |\mathcal{C}_i| > \delta\}_{i=1}^K$$

$$\text{intra}(\mathcal{C}) = \frac{1}{K} \sum_i \frac{m_i}{n_i^2 - n_i}, \quad \text{inter}(\mathcal{C}) = \frac{1}{K(K-1)/2} \sum_{i,j} \frac{m_{ij}}{2n_i n_j}$$



• Experiment ⑤: WSC by design can naturally prune noise/outlier samples.



Findings:

- It took 2 min on a Intel i7-CPU machine to get WSE for 200K images
- WSC is able to rectify incorrent weak annotations by pre-trained classifiers.



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 WSC-annotated images

AU	①			②		
	AlexNet {gt15k wlb10k}	AlexNet {gt15k wsc10k}	AlexNet {gt25k}	DRML {gt25k}	AlexNet {gt25k wsc25k}	DRML {gt25k wsc25k}
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]

AU	③			④		
	F1 {gt25k}	LapSVM {gt25k wlb10k}	TSVM {gt25k wlb25k}	S score {gt25k}	LapSVM {gt25k wlb10k}	TSVM {gt25k wlb25k}
1	19.3	1.2	24.1	66.1	82.3	70.2
4	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	79.1
12	85.1	57.7	87.4	75.5	50.9	80.2
25	85.8	88.9	88.2	72.4	79.4	78.5
26	39.0	5.0	47.0	69.5	83.2	78.4
Avg.	52.2	37.0	56.4	68.2	73.1	75.6

AU	20k		200k		400k		1M	
	wlb	wsc	wlb	wsc	wlb	wsc	wlb	wsc
1	17.6	18.3	17.8	19.3	16.9	[21.3]	17.6	21.2
4	20.3	20.5	19.0	20.4	18.9	21.3	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.