

Label Propagation for Deep Semi-supervised Learning

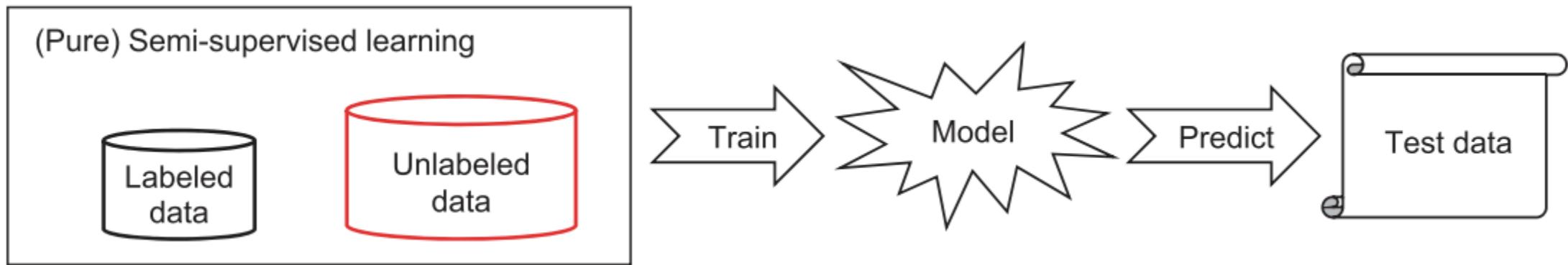
Ahmet Iscen, Giorgos Tolias, Yannis Avrithis & Ondrej Chum

Czech Technical University in Prague & Univ Rennes, Inria, CNRS, IRISA

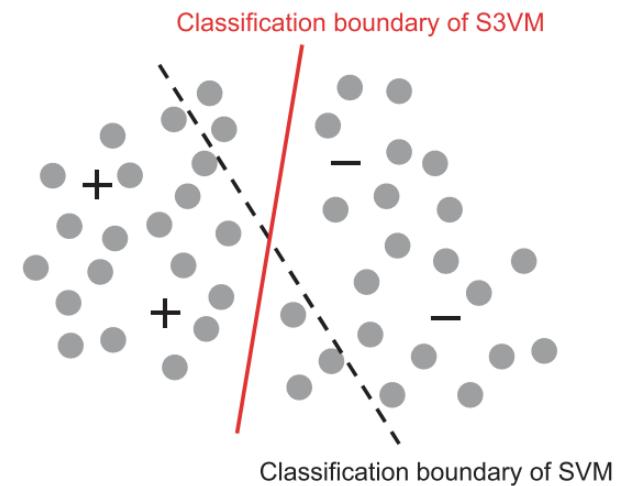
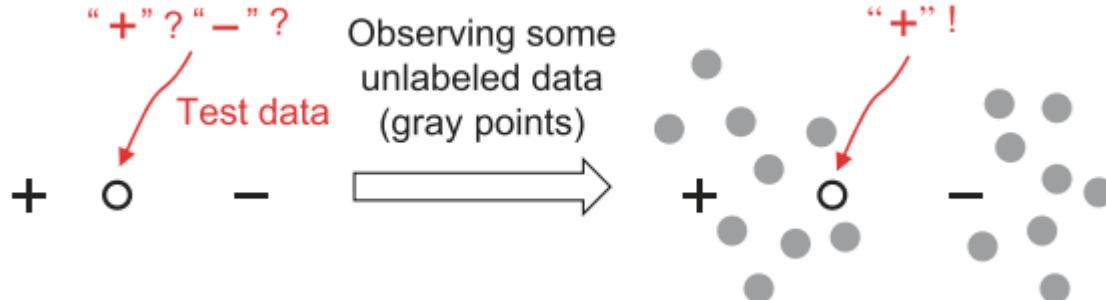
Presented by Xun Xu

Motivations – What is Semi-Supervised Learning (SSL) [1]

- Objective of Semi-Supervised Learning

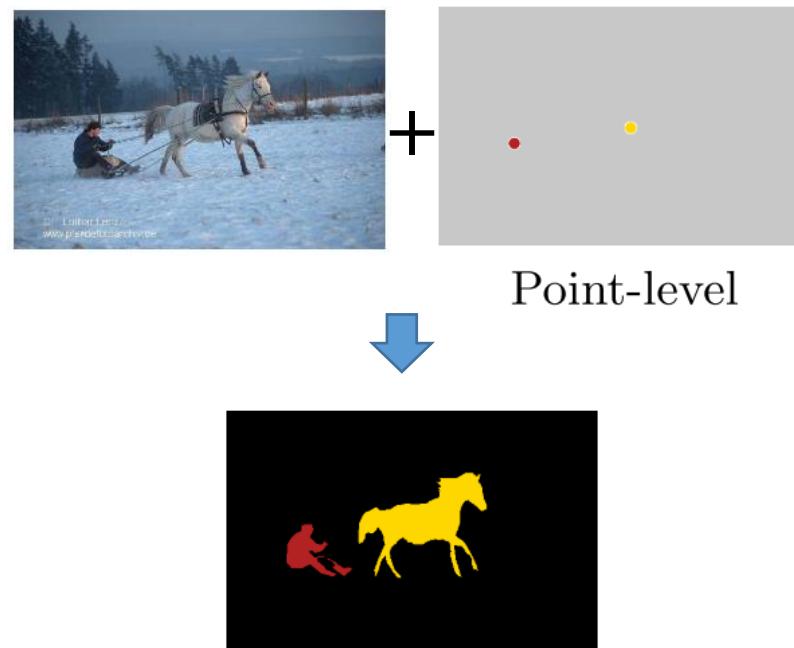
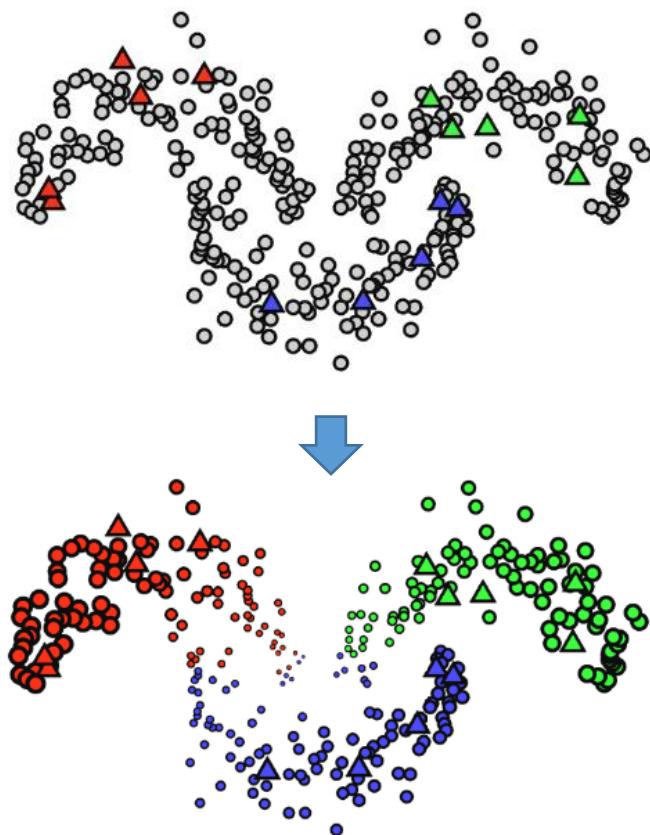


- Why it works?

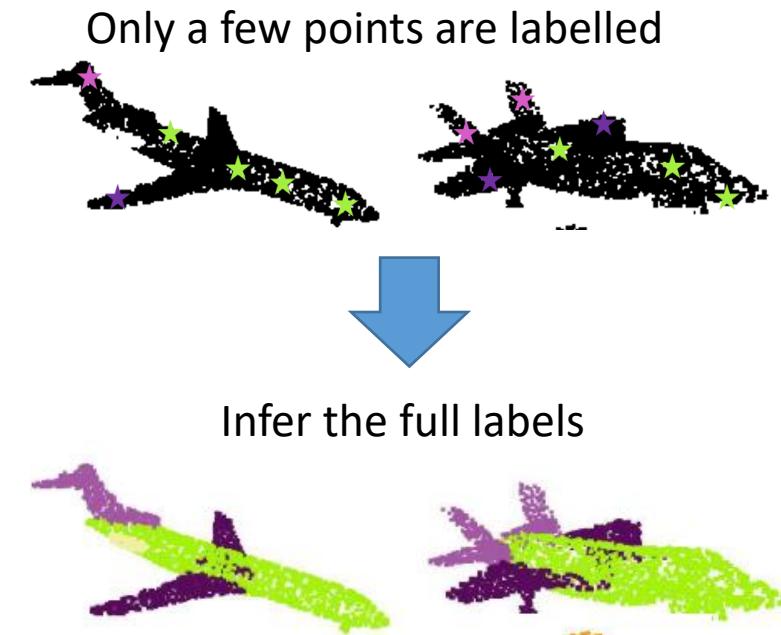


Motivation - Examples

- Examples of SSL on toy example & computer vision data



Point-level



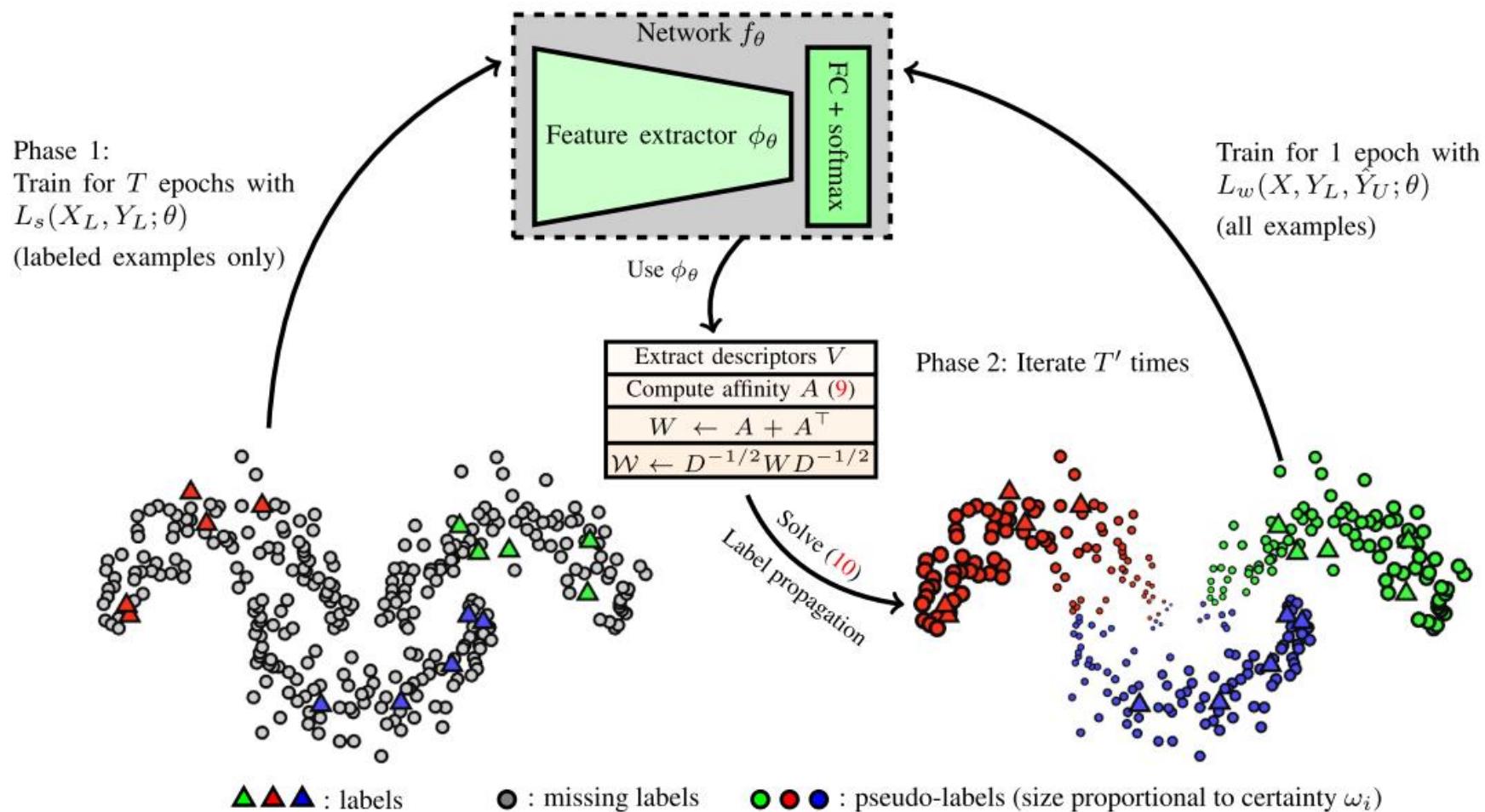
Infer the full labels

Problem Formulation

- Notations
- Input examples: $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_l, \mathbf{x}_{l+1}, \dots, \mathbf{x}_n\}$
- Input labels: $\mathbf{Y}_l = \{y_1, \dots, y_l\}, y_i \in \{1, \dots, C\}$
- Labelled examples: $\mathbf{X}_l = \{\mathbf{x}_i \mid i \in \{1, \dots, l\}\}$
- Unlabelled examples: $\mathbf{X}_u = \{\mathbf{x}_i \mid i \in \{l+1, \dots, n\}\}$
- Labelled/Unlabelled Set: $\mathbf{L} = \{1, \dots, l\} \quad \mathbf{U} = \{l+1, \dots, n\}$
- Adjacency/Weight Matrix: $\mathbf{W} \in \mathbf{R}^{N \times N}$
- Degree Matrix: $\mathbf{D} = diag(\mathbf{d}) \quad d_i = \sum_j w_{ij}$

Overview of Algorithm

- Algorithm overview for self-supervised label propagation



K Nearest Neighbor Graph

- Encoder network

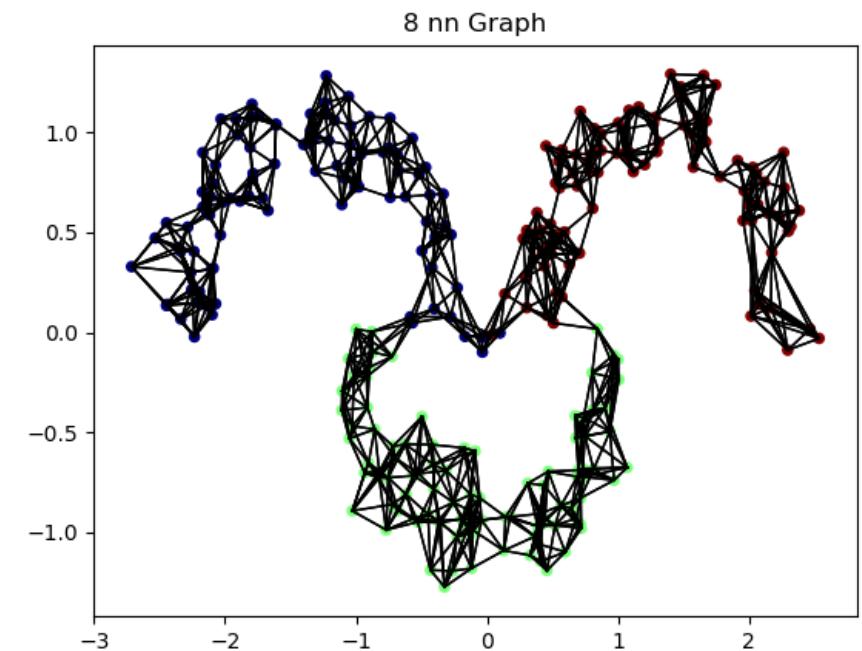
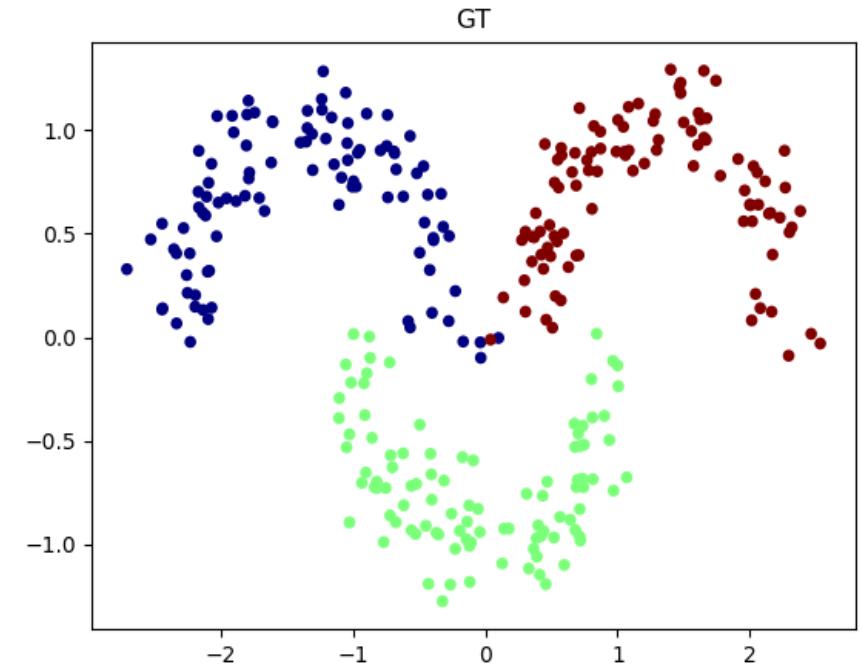
$$\mathbf{v}_i = \phi(\mathbf{x}_i; \boldsymbol{\theta})$$

- Knn Graph

$$a_{ij} = \begin{cases} \exp(-\mathbf{v}_i^T \mathbf{v}_j / \lambda) & \mathbf{v}_j \in NN_k(\mathbf{v}_i) \\ 0 & otherwise \end{cases}$$

- Symmetric Knn Graph

$$\mathbf{W} = \mathbf{A} + \mathbf{A}^T$$



Label Propagation

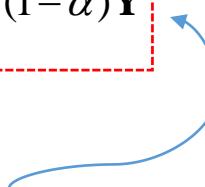
- Normalize Adajency Matrix $\tilde{\mathbf{W}} = \mathbf{D}^{-0.5} \mathbf{W} \mathbf{D}^{-0.5}$
- One-hot Encoded Label Matrix $y_{ij} = \begin{cases} 1 & i \in \{1, \dots, l\} \wedge y_i = j \\ 0 & otherwise \end{cases}, \quad Y \in \{0,1\}^{l+u \times C}$
- Closed-form Label Predictions $\mathbf{Z} = (\mathbf{I} - \alpha \tilde{\mathbf{W}})^{-1} \mathbf{Y} \quad \text{Eq(6) in the paper}$
- Class Prediction $\tilde{y}_i = \arg \max_j z_{ij}$

Label Propagation 15 Years Ago: An Algorithmic Perspective [1,2]

- Label Propagation Algorithm

- Step 1 Compute normalized weight matrix $\tilde{\mathbf{W}} = \mathbf{D}^{-0.5} \mathbf{W} \mathbf{D}^{-0.5}$ Adjacency Matrix \mathbf{W}
- Step 2 Propagate labels according to rule $\mathbf{Z}^{(t+1)} = \alpha \tilde{\mathbf{W}} \mathbf{Z}^t + (1 - \alpha) \mathbf{Y}$ Balance Param. $\alpha \in [0, 1]$
- Step 3 Repeat previous step until convergence Initial Labels \mathbf{Y}
- Step 4 Assign labels $y_i = \arg \max_j \mathbf{Z}_{ij}^\infty$ Propagated Labels \mathbf{Z}

- Closed-Form Solution for Step 2 & Step 3

$$\begin{aligned}\mathbf{Z}^\infty &= \alpha \tilde{\mathbf{W}} \left(\alpha \tilde{\mathbf{W}} \left(\alpha \tilde{\mathbf{W}} (\cdots) + (1 - \alpha) \mathbf{Y} \right) + (1 - \alpha) \mathbf{Y} \right) + (1 - \alpha) \mathbf{Y} \\ &= \alpha^\infty \tilde{\mathbf{W}}^\infty \mathbf{Z}^0 + \sum_{m=0 \dots \infty} \alpha^m \tilde{\mathbf{W}}^m (1 - \alpha) \mathbf{Y} = \boxed{\sum_{m=0 \dots \infty} \alpha^m \tilde{\mathbf{W}}^m (1 - \alpha) \mathbf{Y}} \\ &= (1 - \alpha) (\mathbf{I} - \alpha \tilde{\mathbf{W}})^{-1} \mathbf{Y} \propto (\mathbf{I} - \alpha \tilde{\mathbf{W}})^{-1} \mathbf{Y}\end{aligned}$$


[1] Zhu, Xiaojin, Zoubin Ghahramani, and John D. Lafferty. "Semi-supervised learning using gaussian fields and harmonic functions." *ICML03*

[2] Zhou, Dengyong, et al. "Learning with local and global consistency." *NIPS04*

Proof of Normalized Matrix's Eigenvalue Bounds

- Normalized Weight Matrix

$$\tilde{\mathbf{W}} = \mathbf{D}^{-0.5} \mathbf{W} \mathbf{D}^{-0.5} \text{ so } |eig(\tilde{\mathbf{W}})| \leq 1$$

Proof: [1]

Laplacian Matrix $\mathbf{L} = \mathbf{I} - \tilde{\mathbf{W}}$ $\succ 0$

$$\forall \mathbf{x} \in \mathbf{R}^N \quad \mathbf{x}^T \mathbf{L} \mathbf{x} = \sum_{i=1 \dots N} x_i^2 - \sum_{i,j} \tilde{w}_{ij} x_i x_j = \frac{1}{2} \sum_{i,j} \tilde{w}_{ij} (x_i - x_j)^2 \geq 0$$

Thus the Rayleigh Quotient is upper bounded by 1

$$\mathbf{x}^T \mathbf{L} \mathbf{x} = \mathbf{x}^T (\mathbf{I} - \tilde{\mathbf{W}}) \mathbf{x} = \mathbf{x}^T \mathbf{x} - \mathbf{x}^T \tilde{\mathbf{W}} \mathbf{x} \geq 0 \Rightarrow 1 \geq \frac{\mathbf{x}^T \tilde{\mathbf{W}} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$

The Rayleigh Quotient is lower bounded by -1

$$\mathbf{x}^T (\mathbf{I} + \tilde{\mathbf{W}}) \mathbf{x} = \mathbf{x}^T \mathbf{x} + \mathbf{x}^T \tilde{\mathbf{W}} \mathbf{x} \geq 0 \Rightarrow -1 \leq \frac{\mathbf{x}^T \tilde{\mathbf{W}} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$

The Rayleigh Quotient is bounded by the minimal and maximal eigenvalues of $\tilde{\mathbf{W}}$

Proof of Rayleigh Quotient Bound

- Consider the following problem (minimizing Rayleigh Quotient)

$$\min_{\mathbf{x}} \frac{\mathbf{x}^T \tilde{\mathbf{W}} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \quad s.t. \quad \mathbf{x}^T \mathbf{x} = 1 \Rightarrow \min_{\mathbf{x}} \mathbf{x}^T \tilde{\mathbf{W}} \mathbf{x} \quad s.t. \quad \mathbf{x}^T \mathbf{x} = 1$$

Lagrangian Multiplier $L(\mathbf{x}, \lambda) = \mathbf{x}^T \tilde{\mathbf{W}} \mathbf{x} - \lambda (\mathbf{x}^T \mathbf{x} - 1)$

Derivative of Lagrangian $\frac{\partial L}{\partial \mathbf{x}} = \tilde{\mathbf{W}} \mathbf{x} - \lambda \mathbf{x} = 0 \Rightarrow \tilde{\mathbf{W}} \mathbf{x} = \lambda \mathbf{x}$

- Stationary Point Only at Eigenvectors (same applied to maximizing Rayleigh Quotient)

$$\lambda_{\min} \leq \mathbf{x}^T \tilde{\mathbf{W}} \mathbf{x} \leq \lambda_{\max} \quad s.t. \quad \mathbf{x}^T \mathbf{x} = 1$$

$$\Rightarrow \lambda_{\min} \leq \frac{\mathbf{x}^T \tilde{\mathbf{W}} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \leq \lambda_{\max}$$

Proof of Label Propagation Convergence

- Normalized Weight Matrix

$$\tilde{\mathbf{W}} = \mathbf{D}^{-0.5} \mathbf{W} \mathbf{D}^{-0.5} \text{ so } |eig(\tilde{\mathbf{W}})| \leq 1$$

and $\tilde{\mathbf{W}} = \mathbf{P} \Lambda \mathbf{P}^{-1}$

where $\Lambda = diag(eig(\tilde{\mathbf{W}}))$

- Solve the summation

$$\mathbf{Z}^\infty = \sum_{m=0}^{\infty} \alpha^m \tilde{\mathbf{W}}^m (1-\alpha) \mathbf{Y}$$

$$\begin{aligned} \sum_{m=0}^{\infty} \alpha^m \tilde{\mathbf{W}}^m &= \mathbf{I} + \alpha \mathbf{P} \Lambda \mathbf{P}^{-1} + \alpha^2 \mathbf{P} \Lambda \mathbf{P}^{-1} \mathbf{P} \Lambda \mathbf{P}^{-1} + \alpha^3 \mathbf{P} \Lambda \mathbf{P}^{-1} \mathbf{P} \Lambda \mathbf{P}^{-1} \mathbf{P} \Lambda \mathbf{P}^{-1} + \dots \\ &= \mathbf{I} + \alpha \mathbf{P} \Lambda \mathbf{P}^{-1} + \alpha^2 \mathbf{P} \Lambda^2 \mathbf{P}^{-1} + \alpha^3 \mathbf{P} \Lambda^3 \mathbf{P}^{-1} + \dots \\ &= \mathbf{P} \sum_{m=0}^{\infty} \alpha^m \Lambda^m \mathbf{P}^{-1} \end{aligned}$$

$$= \mathbf{P} \begin{bmatrix} \sum_m (\alpha \lambda_1)^m & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sum_m (\alpha \lambda_N)^m \end{bmatrix} \mathbf{P}^{-1} = \mathbf{P} \begin{bmatrix} \frac{1}{1-\alpha \lambda_1} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \frac{1}{1-\alpha \lambda_N} \end{bmatrix} \mathbf{P}^{-1}$$

$$\begin{aligned} &= \left(\mathbf{P} \begin{bmatrix} 1-\alpha \lambda_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & 1-\alpha \lambda_N \end{bmatrix} \mathbf{P}^{-1} \right)^{-1} = \left(\mathbf{I} - \alpha \mathbf{P} \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \lambda_N \end{bmatrix} \mathbf{P}^{-1} \right)^{-1} \\ &= (\mathbf{I} - \alpha \tilde{\mathbf{W}})^{-1} \end{aligned}$$

Proof of Label Propagation Convergence

- Alternatively, the convergence of label propagation is easily verified as
 - Take the limit on both sides

$$\lim_{t \rightarrow \infty} \mathbf{Z}^{(t+1)} = \lim_{t \rightarrow \infty} \alpha \tilde{\mathbf{W}} \mathbf{Z}^t + (1 - \alpha) \mathbf{Y}$$

$$\Rightarrow \mathbf{Z}^\infty = \alpha \tilde{\mathbf{W}} \mathbf{Z}^\infty + (1 - \alpha) \mathbf{Y}$$

$$\Rightarrow (\mathbf{I} - \alpha \tilde{\mathbf{W}}) \mathbf{Z}^\infty = (1 - \alpha) \mathbf{Y}$$

$$\Rightarrow \mathbf{Z}^\infty = (1 - \alpha) (\mathbf{I} - \alpha \tilde{\mathbf{W}})^{-1} \mathbf{Y}$$

Label Propagation: Optimization Perspective

- The following objective is minimize for the purpose of label propagation

$$\min_{\{\mathbf{z}_i\}} \frac{1}{2} \left(\underbrace{\sum_i \sum_j w_{ij} \left\| \frac{\mathbf{z}_i}{\sqrt{d_i}} - \frac{\mathbf{z}_j}{\sqrt{d_j}} \right\|_2^2}_{\text{Smoothness Constraint}} + \underbrace{\beta \sum_i \|\mathbf{z}_i - \mathbf{y}_i\|_2^2}_{\text{Fitting Constraint}} \right) \quad \text{Eq(8) in the paper}$$

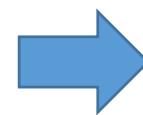
- Smoothness Constraint: encourages labels to be smooth on the manifold
- Fitting Constraint: encourages predictions to be consistent with initial labels

Label Propagation: Optimization Perspective

- Derive the closed-form solution

Notation: $\mathbf{Z} = [\mathbf{z}_1 \ \cdots \ \mathbf{z}_N]$

$$\begin{aligned}
 c &= \frac{1}{2} \sum_i \sum_j w_{ij} \left\| \frac{\mathbf{z}_i}{\sqrt{d_i}} - \frac{\mathbf{z}_j}{\sqrt{d_j}} \right\|_2^2 + \beta \sum_i \|\mathbf{z}_i - \mathbf{y}_i\|_2^2 \\
 &= \frac{1}{2} \sum_i \sum_j \left(w_{ij} \frac{\mathbf{z}_i^T \mathbf{z}_i}{d_i} + w_{ij} \frac{\mathbf{z}_j^T \mathbf{z}_j}{d_j} - 2w_{ij} \frac{\mathbf{z}_i^T \mathbf{z}_j}{\sqrt{d_i d_j}} \right) + \beta \sum_i (\mathbf{z}_i^T \mathbf{z}_i - 2\mathbf{z}_i^T \mathbf{y}_i) + C \\
 &= \sum_i \sum_j w_{ij} \frac{\mathbf{z}_i^T \mathbf{z}_i}{d_i} - \sum_i \sum_j w_{ij} \frac{\mathbf{z}_i^T \mathbf{z}_j}{\sqrt{d_i d_j}} + \beta \sum_i (\mathbf{z}_i^T \mathbf{z}_i - 2\mathbf{z}_i^T \mathbf{y}_i) + C \\
 &= \sum_i \mathbf{z}_i^T \mathbf{z}_i - \sum_i \sum_j \tilde{w}_{ij} \mathbf{z}_i^T \mathbf{z}_j + \beta \sum_i (\mathbf{z}_i^T \mathbf{z}_i - 2\mathbf{z}_i^T \mathbf{y}_i) + C \\
 &= \text{tr}(\mathbf{Z}^T \mathbf{Z}) - \text{tr}(\tilde{\mathbf{W}}^T \mathbf{Z}^T \mathbf{Z}) + \beta \text{tr}(\mathbf{Z}^T \mathbf{Z}) - \beta \text{tr}(\mathbf{Z}^T \mathbf{Y}) + C
 \end{aligned}$$



$$\begin{aligned}
 \frac{\partial c}{\partial \mathbf{Z}} &= \mathbf{Z} - \tilde{\mathbf{W}}\mathbf{Z} + \beta\mathbf{Z} - \beta\mathbf{Y} = 0 \\
 \Rightarrow \left(\mathbf{I} - \frac{1}{1+\beta} \tilde{\mathbf{W}} \right) \mathbf{Z} &= \frac{\beta}{1+\beta} \mathbf{Y} \\
 \Rightarrow \mathbf{Z} &= \beta \alpha \left(\mathbf{I} - \alpha \tilde{\mathbf{W}} \right)^{-1} \mathbf{Y} \quad s.t. \alpha = \frac{1}{1+\beta}
 \end{aligned}$$

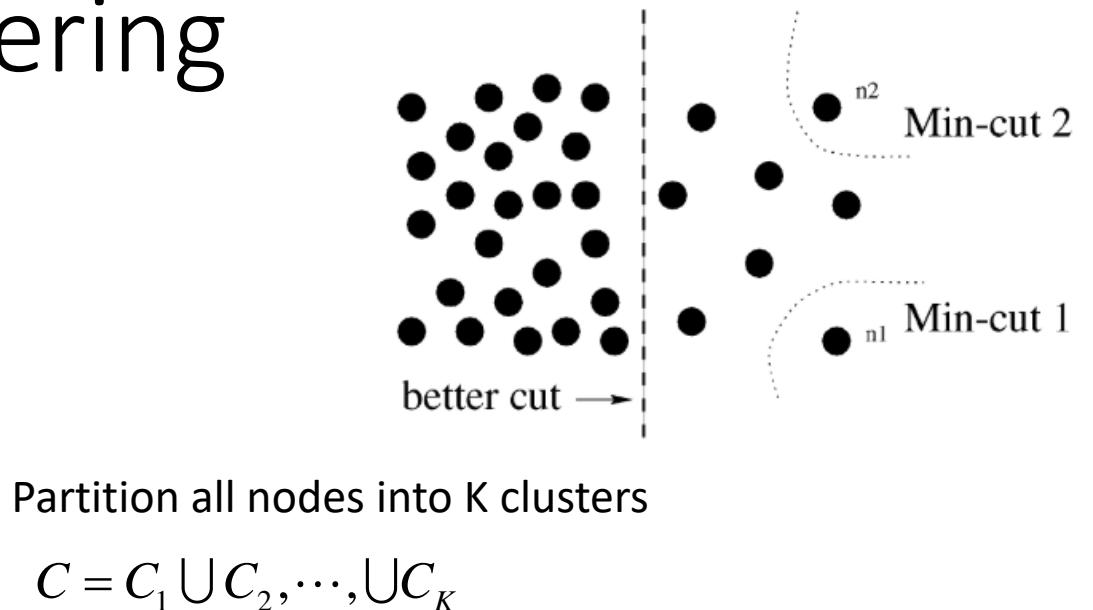
Relation to Spectral Clustering

- Normalized Cut (Ncut)[1]

$$Ncut(C_1, \dots, C_K) = \sum_{k=1 \dots K} \frac{cut(C_k, \bar{C}_k)}{vol(C_k)} = \frac{1}{2} \sum_{k=1 \dots K} \frac{W(C_k, \bar{C}_k)}{vol(C_k)}$$

Denote $f_i \in \{0,1\}^K$ as the partition index for node i

$$W(C_k, \bar{C}_k) = \sum_{i \in C_k} \sum_{j \in \bar{C}_k} w_{ij} \|f_i - f_j\|_2^2 = \sum_{i \in C_k} \sum_{j \in C} w_{ij} \|f_i - f_j\|_2^2$$



$$vol(C_k) = \sum_{i \in C_k} \sum_{j \in C} w_{ij} = \sum_{i \in C_k} d_i$$

Relation to Spectral Clustering - NCut

$$\begin{aligned}
 Ncut(C_1, \dots, C_K) &= \frac{1}{2} \sum_{k=1 \dots K} \frac{\sum_{i \in C_k} \sum_{j \in C} w_{ij} \|f_i - f_j\|_2^2}{\sum_{i \in C_k} d_i} \\
 &= \frac{1}{2} \sum_{k=1 \dots K} \sum_{i \in C_k} \sum_{j \in C} w_{ij} \|\tilde{f}_i - \tilde{f}_j\|_2^2 = \frac{1}{2} \sum_{i \in C} \sum_{j \in C} w_{ij} \|\tilde{f}_i - \tilde{f}_j\|_2^2 \\
 &= \frac{1}{2} \sum_{i \in C} \sum_{j \in C} w_{ij} \left\| \frac{\mathbf{z}_i}{\sqrt{d_i}} - \frac{\mathbf{z}_j}{\sqrt{d_j}} \right\|_2^2
 \end{aligned}$$

The Final Objective

$$\begin{aligned}
 &\min_{\{f_i\}} Ncut(C_1, \dots, C_K) \\
 &\Rightarrow \min_{\{\mathbf{z}_i\}} \frac{1}{2} \sum_{i \in C} \sum_{j \in C} w_{ij} \left\| \frac{\mathbf{z}_i}{\sqrt{d_i}} - \frac{\mathbf{z}_j}{\sqrt{d_j}} \right\|_2^2 \quad s.t. \quad \mathbf{Z}^T \mathbf{Z} = \mathbf{I}
 \end{aligned}$$



$$\tilde{f}_i = \frac{\mathbf{f}_i}{\sqrt{\sum_{i \in C_k} d_i}} \quad s.t. \quad i \in C_k$$

$$\tilde{\mathbf{F}} = \left[\underbrace{\tilde{f}_1 \dots \tilde{f}_{|C_1|}}_{|C_1|}, \underbrace{\tilde{f}_{|C_1|+1} \dots \tilde{f}_{|C_1|+|C_2|}}_{|C_2|}, \dots \underbrace{\tilde{f}_{|C_1|+\dots+|C_{K-1}|} \dots \tilde{f}_N}_{|C_K|} \right]^T \quad \tilde{\mathbf{F}}^T \mathbf{D} \tilde{\mathbf{F}} = \mathbf{I}$$

Let $\mathbf{Z} = \mathbf{D}^{0.5} \tilde{\mathbf{F}}$ We have $\mathbf{Z}^T \mathbf{Z} = \mathbf{I}$ and $\tilde{f}_i = \frac{\mathbf{z}_i}{\sqrt{d_i}}$

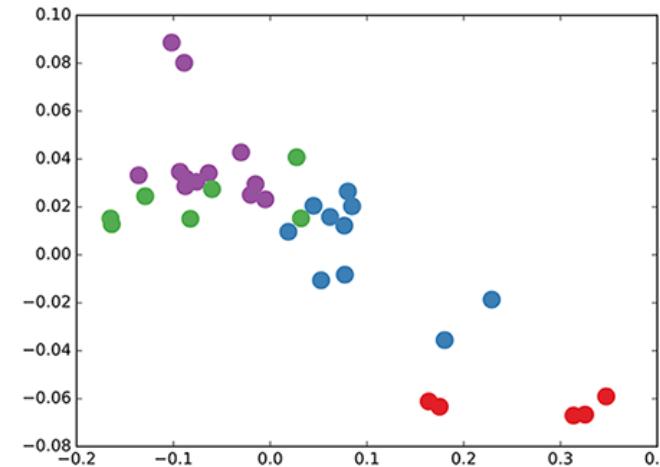
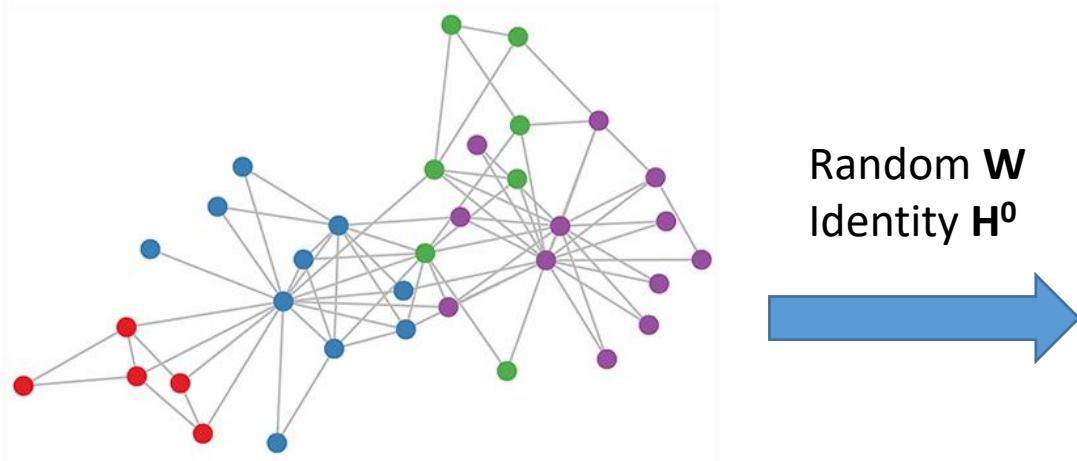
Solve by relaxing $\mathbf{Z} \in \mathbf{R}^{N \times K}$

Relation to GCN

- The Propagation Rule for GCN [1]

$$\mathbf{H}^{(l+1)} = \sigma\left(\tilde{\mathbf{D}}^{-0.5}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-0.5}\mathbf{H}^{(l)}\mathbf{W}^{(l)}\right) \quad \mathbf{D} = \text{diag}\left(\sum_j \tilde{a}_{ij}\right), \quad \tilde{\mathbf{A}} = \mathbf{I} + \mathbf{A}$$

- Explain the Success of Random Initialization [1]



Efficient Computation

- Large-Scale Label Propagation

$$\mathbf{Z} = (\mathbf{I} - \alpha \tilde{\mathbf{W}})^{-1} \mathbf{Y} \quad \text{The inverse is up to } O(n^3)$$

- Iterative Method for $(\mathbf{I} - \alpha \tilde{\mathbf{W}}) \mathbf{Z} = \mathbf{Y}$
- Denote $\mathbf{A} = \mathbf{I} - \alpha \tilde{\mathbf{W}}$, $\mathbf{b} = \mathbf{y}_i$
- Solve the unconstraint quadratic problem $\min_{\mathbf{z}_i} \mathbf{z}_i^T \mathbf{A} \mathbf{z}_i - \mathbf{b} \mathbf{z}_i$
- Unique solution exists $\mathbf{A} \succ 0$
- Gradient Descent $\mathbf{z}_{(t+1)} = \mathbf{z}_{(t)} - \alpha_{(t)} \mathbf{d}_{(t)}$, $s.t. \mathbf{d}_{(t)} = \mathbf{A} \mathbf{z}_{(t)} - \mathbf{b}$
- Conjugate Descent [1]: How to find $\alpha_{(t)}$ and $\mathbf{d}_{(t)}$ to speed up convergence

Self-Supervised Training

- Calculate Weight for Each Sample

$$\omega_i = 1 - \frac{H(\hat{\mathbf{z}}_i)}{\log(C)} \quad \hat{z}_{ij} = z_{ij} / \sum_k z_{ik}$$

- Calculate Weight for Each Class

$$\zeta_j = (\|\mathbf{L}_j\| + \|\mathbf{U}_j\|)^{-1}$$

- The Final Loss

$$Lw(\mathbf{X}, \mathbf{Y}_l, \mathbf{Y}_u; \boldsymbol{\theta}) = \sum_{i=1 \dots l} \zeta_{y_i} l_s(f(\mathbf{x}_i; \boldsymbol{\theta}), \mathbf{y}_i) + \sum_{i=l+1 \dots n} \omega_i \zeta_{\hat{y}_i} l_s(f(\mathbf{x}_i; \boldsymbol{\theta}), \hat{\mathbf{y}}_i)$$

Summary of Algorithm

Algorithm 1 Label propagation for deep SSL

```
1: procedure LPDSSL(Training examples  $X$ , labels  $Y_L$ )
2:    $\theta \leftarrow$  initialize randomly
3:   for epoch  $\in [1, \dots, T]$  do
4:      $\theta \leftarrow \text{OPTIMIZE}(L_s(X_L, Y_L; \theta))$             $\triangleright$  mini-batch optimization
5:   end for
6:   for epoch  $\in [1, \dots, T']$  do
7:     for  $i \in \{1, \dots, n\}$  do  $\mathbf{v}_i \leftarrow \phi_\theta(x_i)$             $\triangleright$  extract descriptors
8:     for  $(i, j) \in \{1, \dots, n\}^2$  do  $a_{ij} \leftarrow$  affinity values (9)
9:      $W \leftarrow A + A^\top$                                  $\triangleright$  symmetric affinity
10:     $W \leftarrow D^{-1/2} W D^{-1/2}$             $\triangleright$  symmetrically normalized affinity
11:     $Z \leftarrow$  solve (10) with CG                       $\triangleright$  diffusion
12:    for  $(i, j) \in U \times C$  do  $\hat{z}_{ij} \leftarrow z_{ij} / \sum_k z_{ik}$             $\triangleright$  normalize  $Z$ 
13:    for  $i \in U$  do  $\hat{y}_i \leftarrow \arg \max_j \hat{z}_{ij}$             $\triangleright$  pseudo-label
14:    for  $i \in U$  do  $\omega_i \leftarrow$  certainty of  $\hat{y}_i$  (11)       $\triangleright$  pseudo-label weight
15:    for  $j \in C$  do  $\zeta_j \leftarrow (|L_j| + |U_j|)^{-1}$         $\triangleright$  class weight/balancing
16:     $\theta \leftarrow \text{OPTIMIZE}(L_w(X, Y_L, \hat{Y}_U; \theta))$             $\triangleright$  mini-batch optimization
17:  end for
18: end procedure
```

Experiments

- Dataset

Dataset	#Classes (tr/te)	#Samples (tr/te)	#Labelled Samples per Class	Image Res.	Network
CIFAR-10	10/10	50k/10k	50/100/200/400	32*32	'13layer' CNN
CIFAR-100	100/100	50k/10k	40/100	32*32	'13layer' CNN
Mini- ImageNet	100/100	50k/10k	40/100	84*84	Resnet-18

- Some Examples



Experiment - Comparisons

- Virtual Adversarial Training [1]
- Adversarial Training with full supervision [2]

$$\min_{\theta} \sum_i D(f(\mathbf{X}_i + \hat{\eta}; \theta), \mathbf{y}_i)$$

D(·, ·) is a divergence and \mathbf{y}_i are known labels
where $\hat{\eta} = \arg \max_{\eta; \|\eta\| \leq \varepsilon} D(f(\mathbf{X}_i + \eta; \theta), \mathbf{y}_i)$

- Virtual Adversarial Training with semi-supervision [1]

$$\min_{\theta} \sum_{i \in \mathbf{L}} D(f(\mathbf{X}_i + \hat{\eta}; \theta), \mathbf{y}_i) + \sum_{i \in \mathbf{U}} D(f(\mathbf{X}_i + \hat{\eta}; \theta^{(t)}), f(\mathbf{X}_i; \theta^{(t-1)}))$$

where $\hat{\eta} = \begin{cases} \arg \max_{\eta; \|\eta\| \leq \varepsilon} D(f(\mathbf{X}_i + \eta; \theta), f(\mathbf{X}_i; \theta^{(t-1)})) & i \in \mathbf{U} \\ \arg \max_{\eta; \|\eta\| \leq \varepsilon} D(f(\mathbf{X}_i + \eta; \theta), \mathbf{y}_i) & i \in \mathbf{L} \end{cases}$

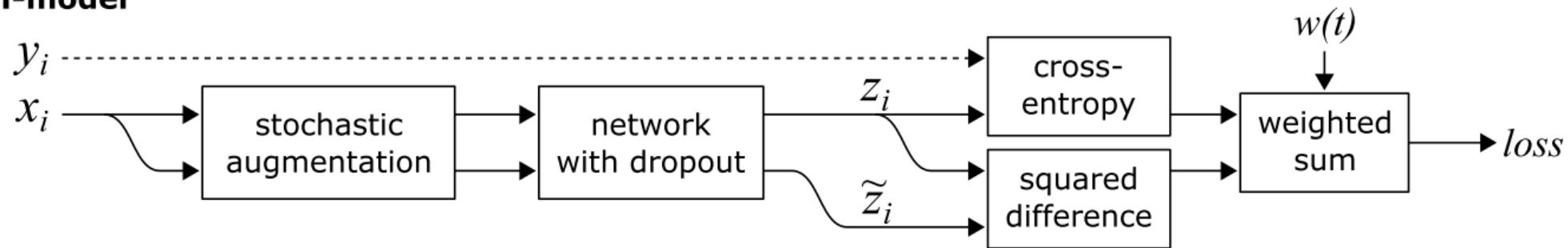
[1] Miyato, Takeru, et al. "Virtual adversarial training: a regularization method for supervised and semi-supervised learning." *TPAMI*. 2018

[2] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." *ICLR*, 2015.

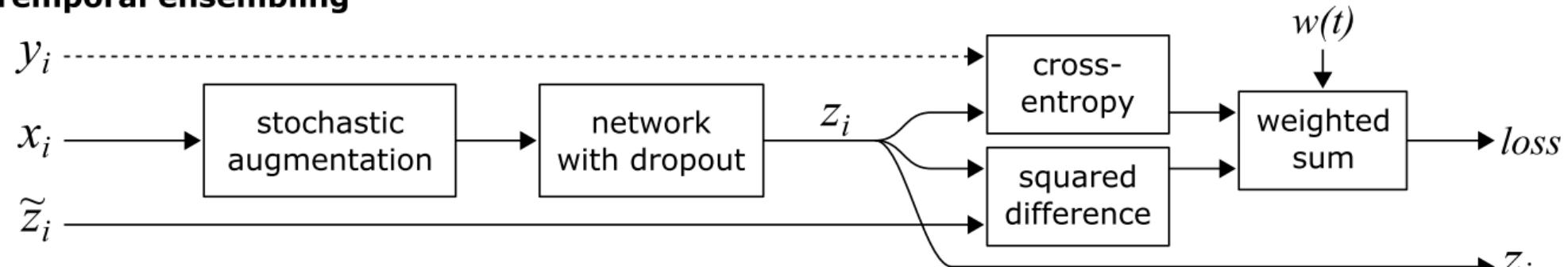
Experiment - Comparisons

- Temporal Ensemble & Π -model [1]

Π -model



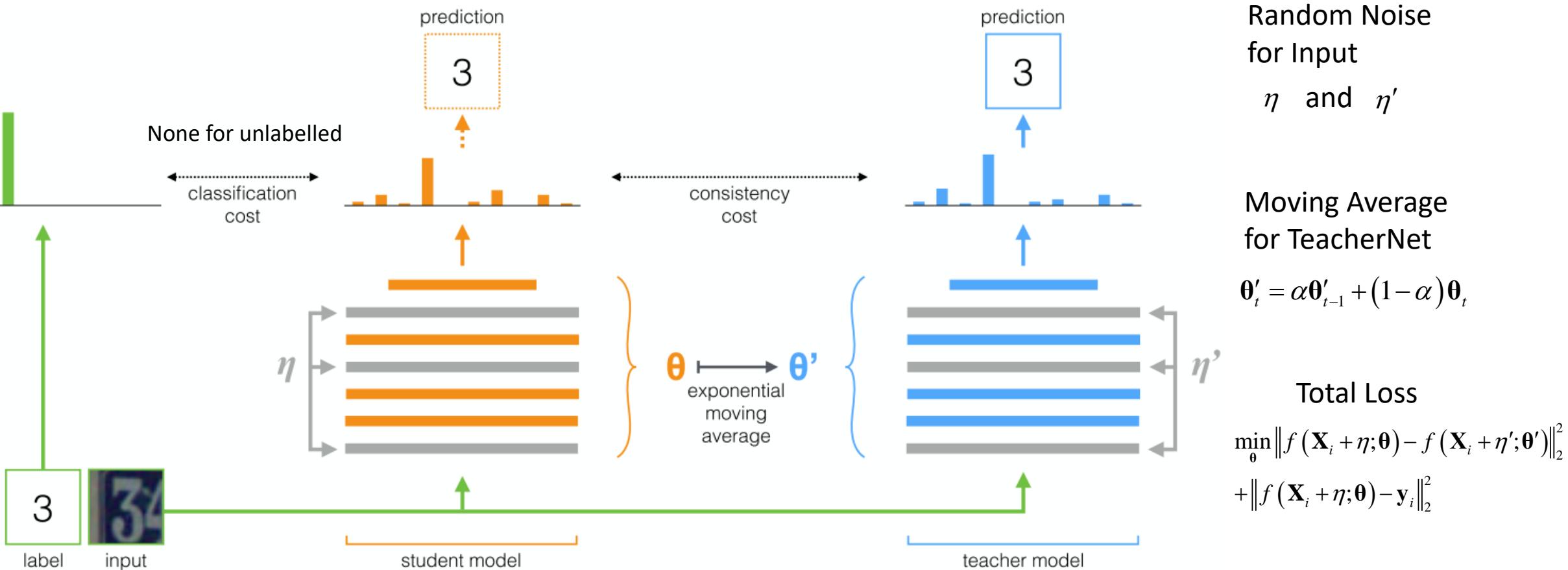
Temporal ensembling



$$\text{Moving Average} \quad \tilde{z}_i^{(t)} = \alpha \tilde{z}_i^{(t-1)} + \frac{(1-\alpha) z_i^{(t-1)}}{1-\alpha^t}$$

Experiment - Comparisons

- Mean Teacher (MT) [1]:



[1] Tarvainen, Antti, and Harri Valpola. "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results." *NIPS*. 2017.

Experiment - Comparisons

- TDCNN [2]:

Training on Labelled Samples

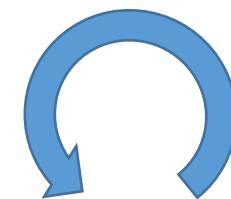
$$\min_{\boldsymbol{\theta}^{(0)}} \sum_{i \in L} CE\left(f\left(\mathbf{X}_i; \boldsymbol{\theta}^{(0)}\right), \mathbf{y}_i\right) \rightarrow \min_{\boldsymbol{\theta}_t} \sum_{i \in U} r_i CE\left(f\left(\mathbf{X}_i; \boldsymbol{\theta}^{(t)}\right), f\left(\mathbf{X}_i; \boldsymbol{\theta}^{(t-1)}\right)\right) + \sum_{i \in L} CE\left(f\left(\mathbf{X}_i; \boldsymbol{\theta}^{(t)}\right), \mathbf{y}_i\right) + L_{reg}$$

Per-sample weight depends on density

$$d_i = \sum_{j \in NNk(f_i)} \|f_i - f_j\|_2^2$$

$$r_i = 1 - \frac{d_i}{d_{max}}$$

Iterative Training



Min-Max Feature Regularization

- Samples of same class to be close
- Samples of different classes to be faraway, above a threshold

$$L_{reg} = r_i r_j \left(\mathbf{1}(y_i = y_j) \|f_i - f_j\|_2^2 - (1 - \mathbf{1}(y_i = y_j)) \min(0, \|f_i - f_j\|_2^2 - \tau) \right)$$

Experiment - Quantitative Evaluation on All 3 Datasets

Dataset	CIFAR-100		Mini-ImageNet- <i>top1</i>		Mini-ImageNet- <i>top5</i>	
	4000	10000	4000	10000	4000	10000
Nb. labeled images	4000	10000	4000	10000	4000	10000
Fully supervised	55.43 ± 0.11	40.67 ± 0.49	74.78 ± 0.33	60.25 ± 0.29	53.07 ± 0.68	38.28 ± 0.38
Ours	46.20 ± 0.76	38.43 ± 1.88	70.29 ± 0.81	57.58 ± 1.47	47.58 ± 0.94	36.14 ± 2.19
MT [38]	45.36 ± 0.49	36.08 ± 0.51	72.51 ± 0.22	57.55 ± 1.11	49.35 ± 0.22	32.51 ± 1.31
MT + Ours	43.73 ± 0.20	35.92 ± 0.47	72.78 ± 0.15	57.35 ± 1.66	50.52 ± 0.39	31.99 ± 0.55

Table 3. Performance comparison on CIFAR-100 and Mini-ImageNet with 4k and 10k labeled images. Error rate is reported. “13-layer” network is used for CIFAR-100 and Resnet-18 is used for Mini-ImageNet. All methods are reproduced by us.

Experiment – Comparison on CIFAR-10

Dataset	CIFAR-10			
	500	1000	2000	4000
Nb. labeled images				
Fully supervised	49.08 ± 0.83	40.03 ± 1.11	29.58 ± 0.93	21.63 ± 0.38
TDCNN [36] [†]	-	32.67 ± 1.93	22.99 ± 0.79	16.17 ± 0.37
Network prediction (1) + weights	35.17 ± 2.46	23.79 ± 1.31	16.64 ± 0.48	13.21 ± 0.61
Ours: Diffusion prediction (7) + weights	32.40 ± 1.80	22.02 ± 0.88	15.66 ± 0.35	12.69 ± 0.29
VAT [26] [†]	-	-	-	11.36
Π model [23] [†]	-	-	-	12.36 ± 0.31
Temporal Ensemble [23] [†]	-	-	-	12.16 ± 0.24
MT [38] [†]	-	27.36 ± 1.30	15.73 ± 0.31	12.31 ± 0.28
MT [38]	27.45 ± 2.64	19.04 ± 0.51	14.35 ± 0.31	11.41 ± 0.25
MT + Ours	24.02 ± 2.44	16.93 ± 0.70	13.22 ± 0.29	10.61 ± 0.28

Table 2. Comparison with the state of the art on CIFAR-10. Error rate is reported. “13-layer” network is used. The top part of the table corresponds to training with pseudo-labels, while the bottom part of the table includes methods that are complementary to ours, as shown by the combination of our method with MT. [†] denotes scores reported in prior work.

Experiments - Ablation Studies

Contribution
of weighted
learning

Pseudo-labeling	ω_i	ζ_j	CIFAR-10
Diffusion (7)		✓	36.53 \pm 1.42
	✓		36.17 \pm 1.98
	✓	✓	33.32 \pm 1.53
GTG [8]	✓	✓	32.40 \pm 1.80
Network (1)	✓	✓	35.20 \pm 2.23
	✓	✓	35.17 \pm 2.46

Table 1. Impact of weights ω_i , class weights ζ_j , and pseudo-labeling by diffusion prediction (7) or network prediction (1). Error rate is reported on CIFAR-10 with 500 labels.

Evolution
of
Weight

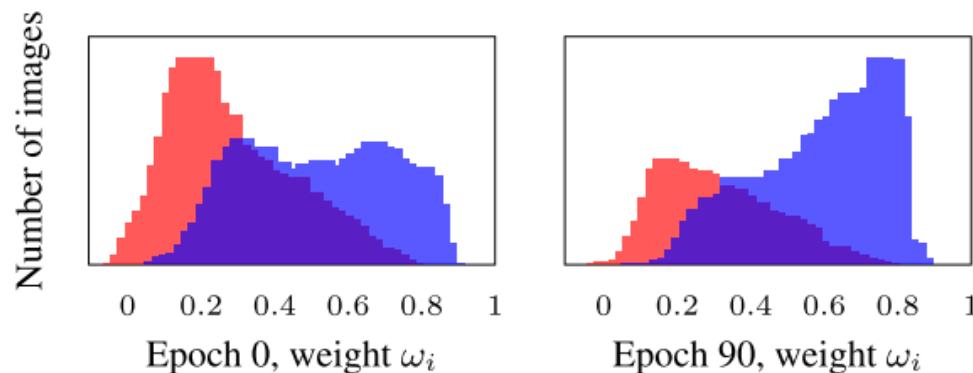


Figure 5. Distribution of weights ω_i for unlabeled images at epoch 0 (left) and epoch 90 (right) during the training of CIFAR-10 with 500 labels. Correct pseudo-labels according to ground-truth are shown in blue and incorrect in red.

Contribution
of Label
Propagation

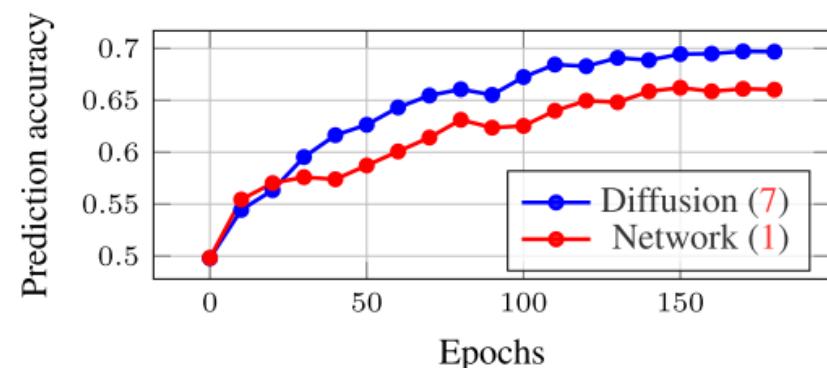


Figure 4. Accuracy of predicted pseudo-labels according to ground-truth on CIFAR-10 with 500 labeled images. Diffusion predictions (7) are compared against network predictions (1).

Varying #
Labelled
Data

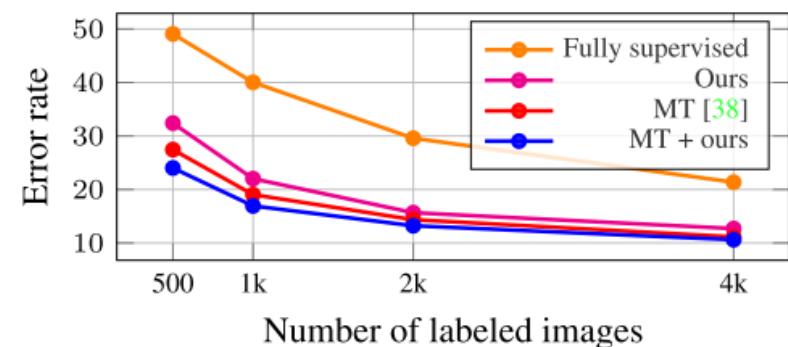
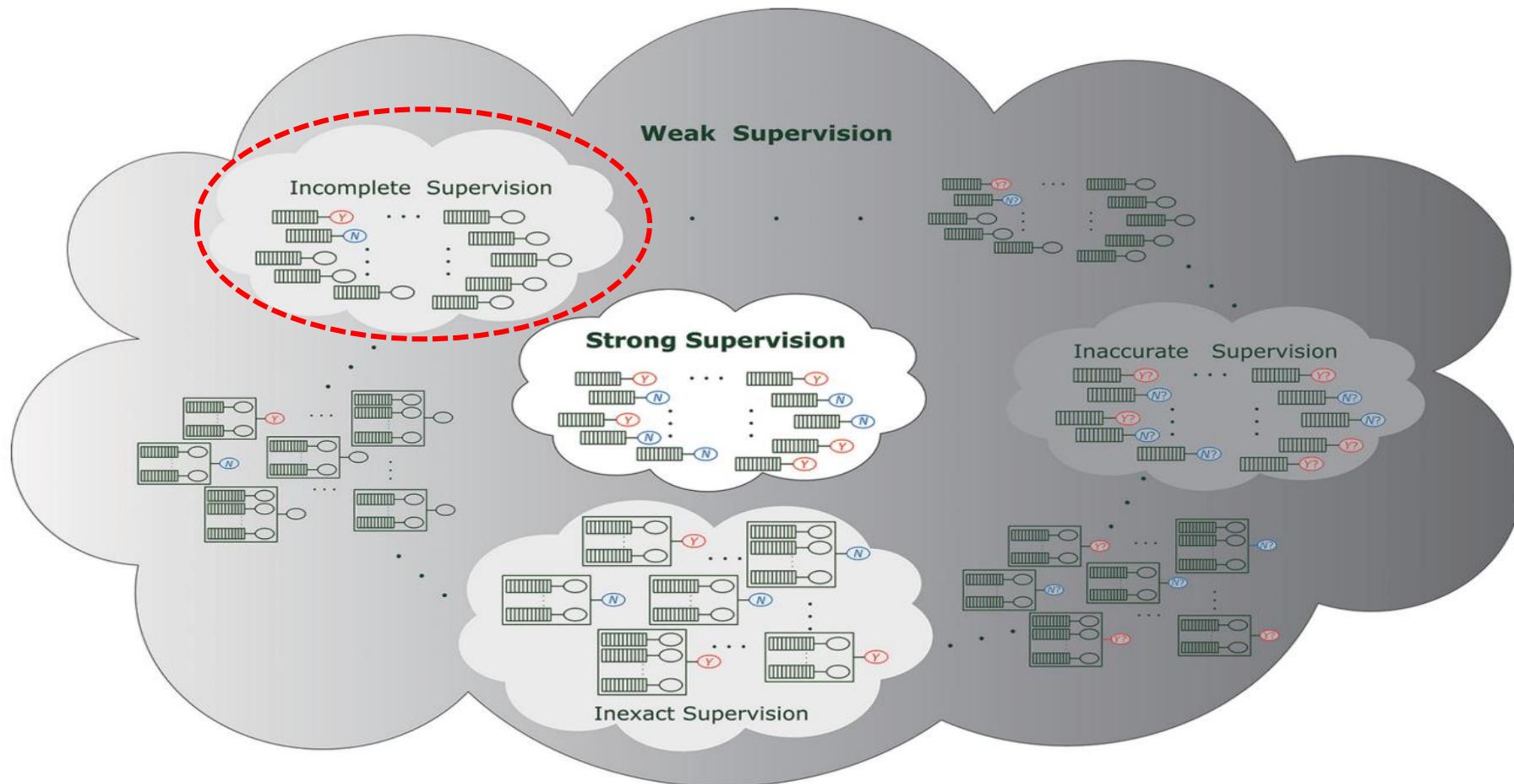


Figure 6. Error rate versus number of labeled images on CIFAR-10 using different methods.

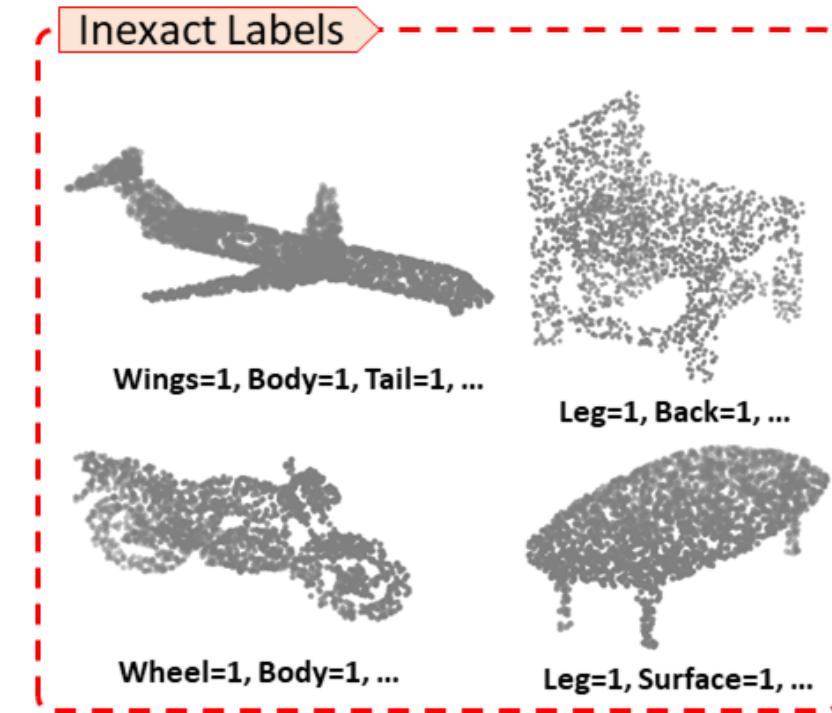
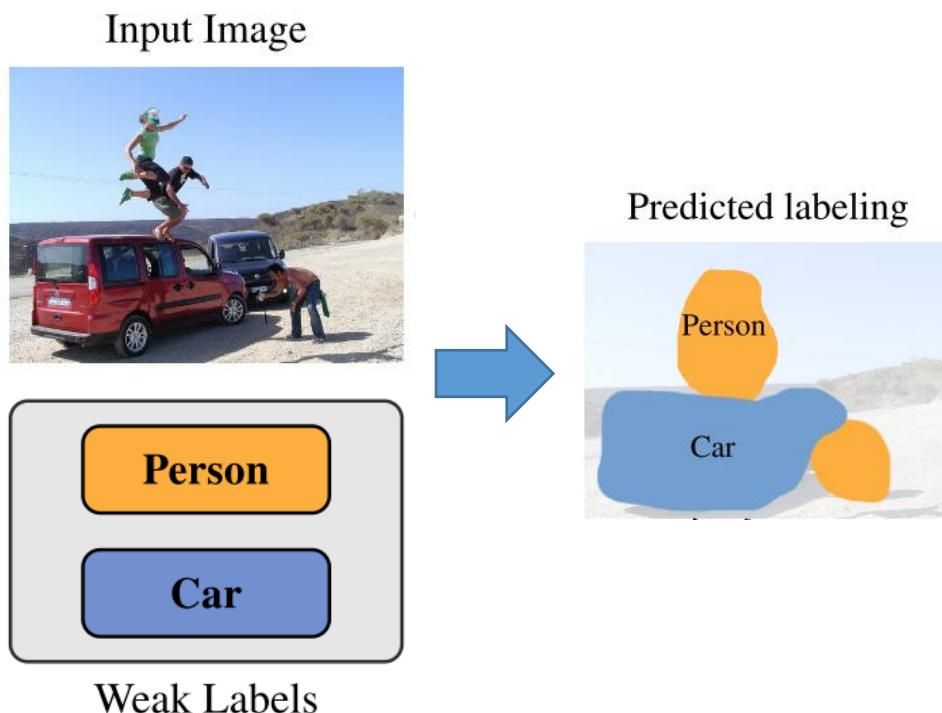
Future Extension

- Semi-Supervised Learning in the Spectrum of Weak Supervision [1]



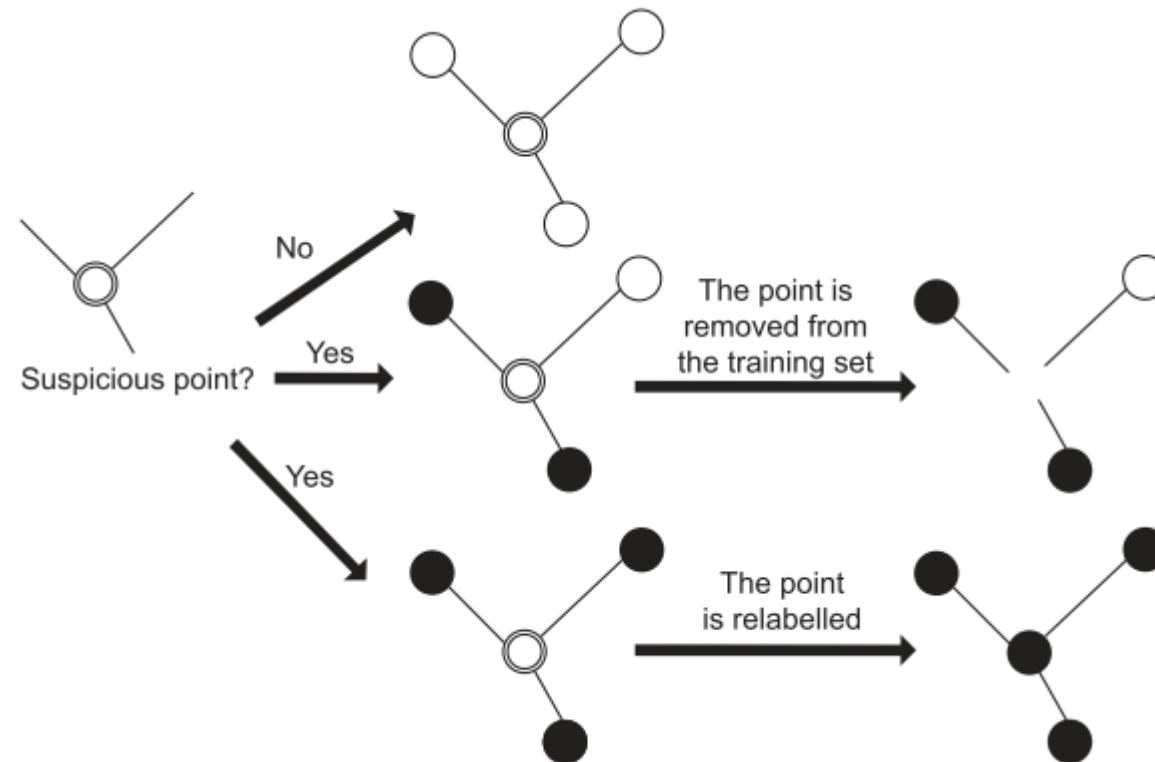
Future Extension

- Inexact Supervision (Weakly Supervised Learning)
- Full annotation is too expensive and is not scalable, e.g. image segmentation & point cloud segmentation



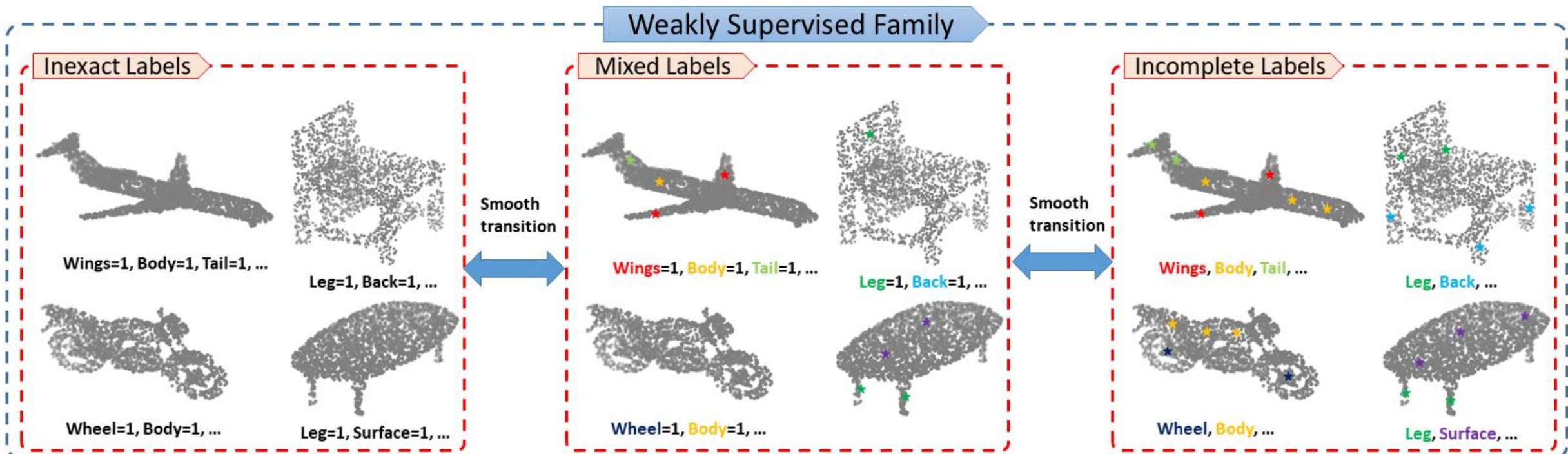
Future Extension

- Inaccurate Supervision (Some Labels are Wrong)
 - Annotation is challenging or annotators are intentional, e.g. per-pixel/point annotation, robot annotation & adversarial annotations



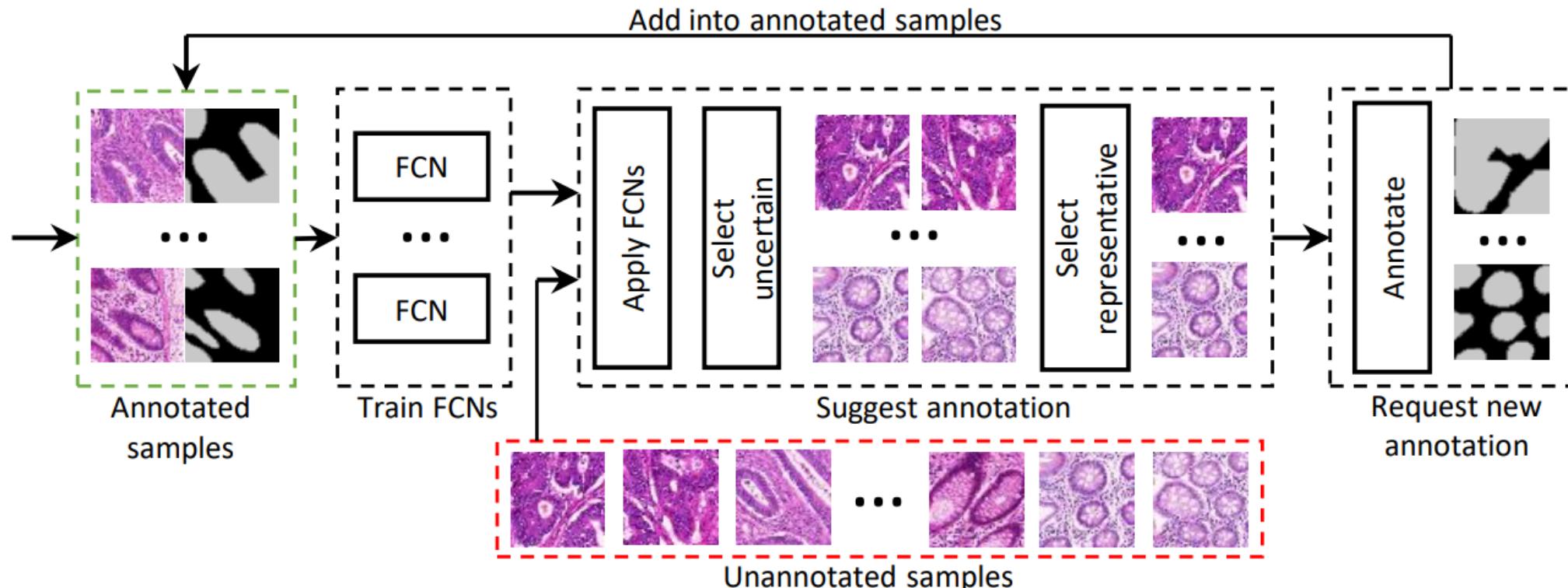
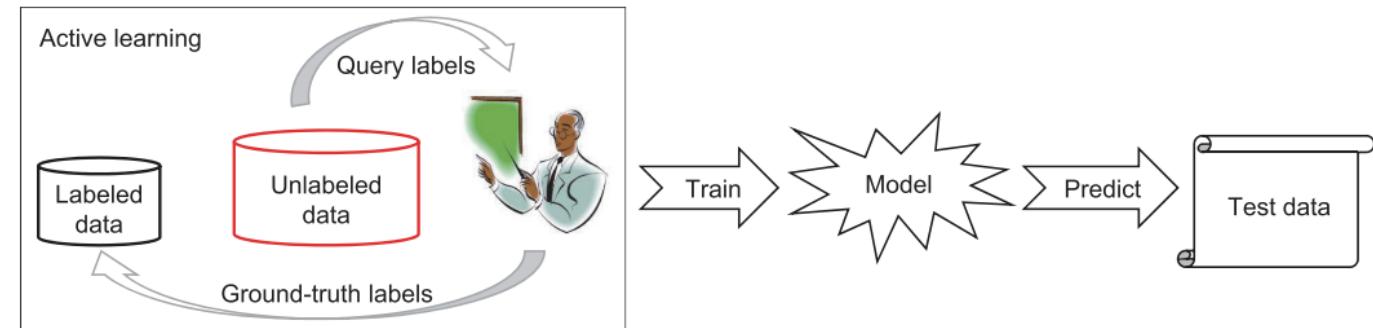
Future Extension

- At the Intersection of Incomplete and Inexact Labels



Future Extension

- Active Learning



[1] Yang, Lin, et al. "Suggestive annotation: A deep active learning framework for biomedical image segmentation." *International conference on medical image computing and computer-assisted intervention*. Springer, Cham, 2017.

Conclusion

- Exploiting unlabelled data for learning tasks is promising
- Relation to traditional graph theorem is revitalizing
- Various forms of weak supervision and/or less costly annotation could be explored