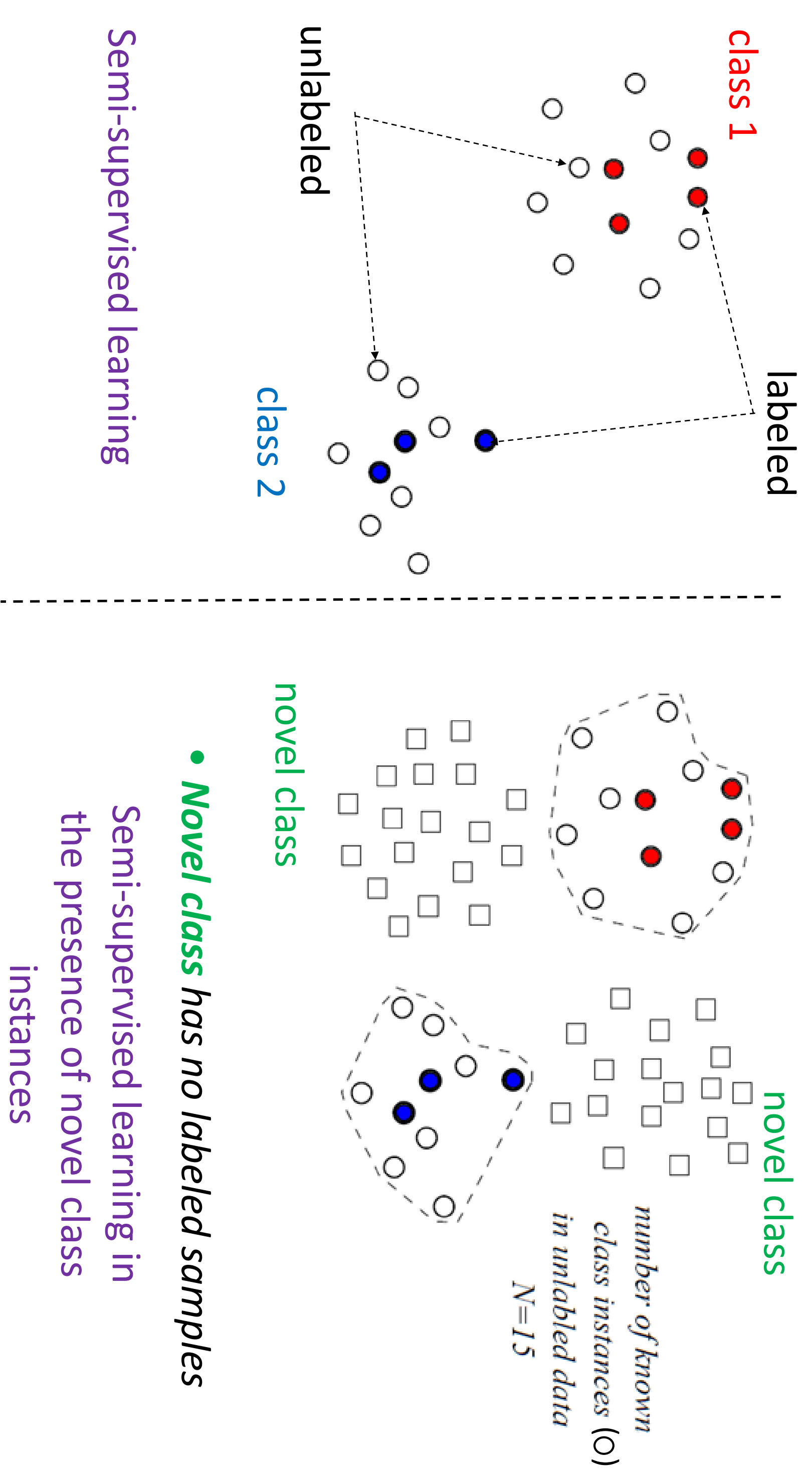
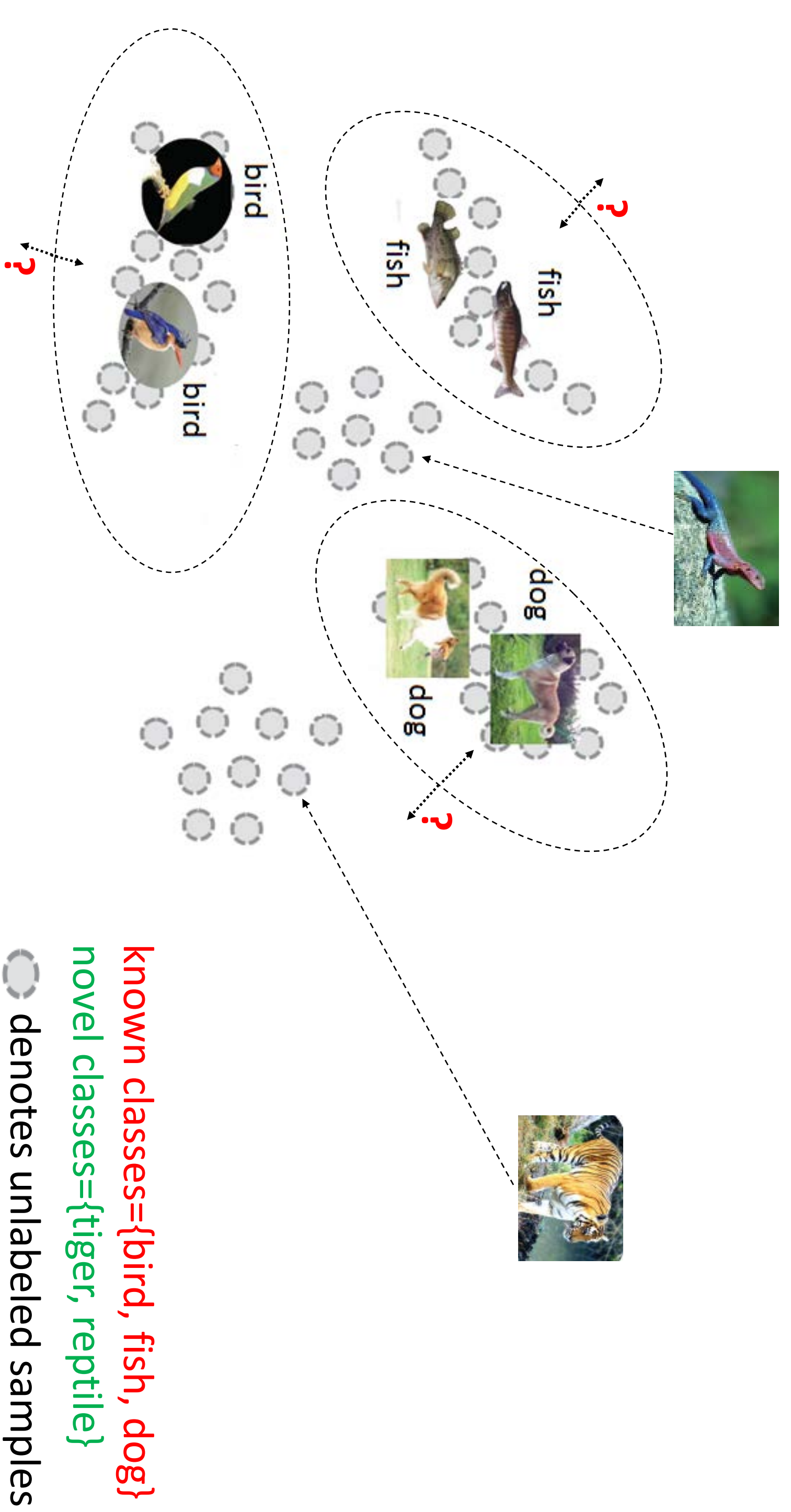


1 Semi-supervised learning with novel class instances



2 Examples in computer vision



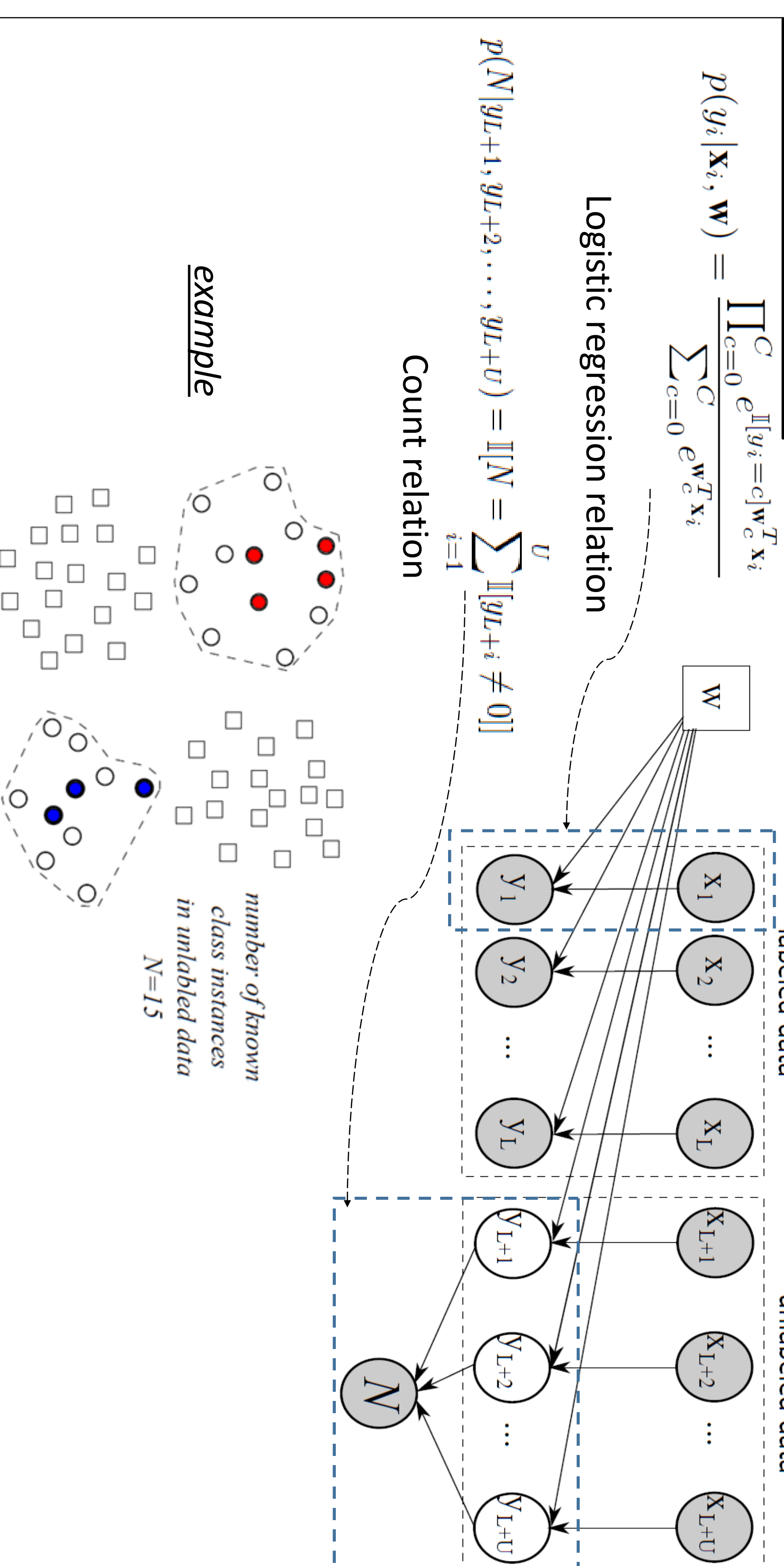
3 Problem formulation

Training data: Labeled data $\{x_i, y_i\}_{i=1}^L$ and unlabeled data $\{x_i\}_{i=1}^U$

1. Labeled data: $x_i \in \mathbb{X} \subseteq \mathbb{R}^d$, $y_i \in \{1, \dots, C\}$ which are known classes.
2. Unlabeled data: $x_{L+i} \in \mathbb{R}^d$, $y_{L+i} \in \{0, 1, \dots, C\}$ which are both known classes and novel class (denoted as 0).
3. Number of novel class instances: N .

Goal: Find a classifier that maps an instance in \mathbb{X} into a label in $\mathbb{Y}=\{0,1, \dots, C\}$

4 Graphical model



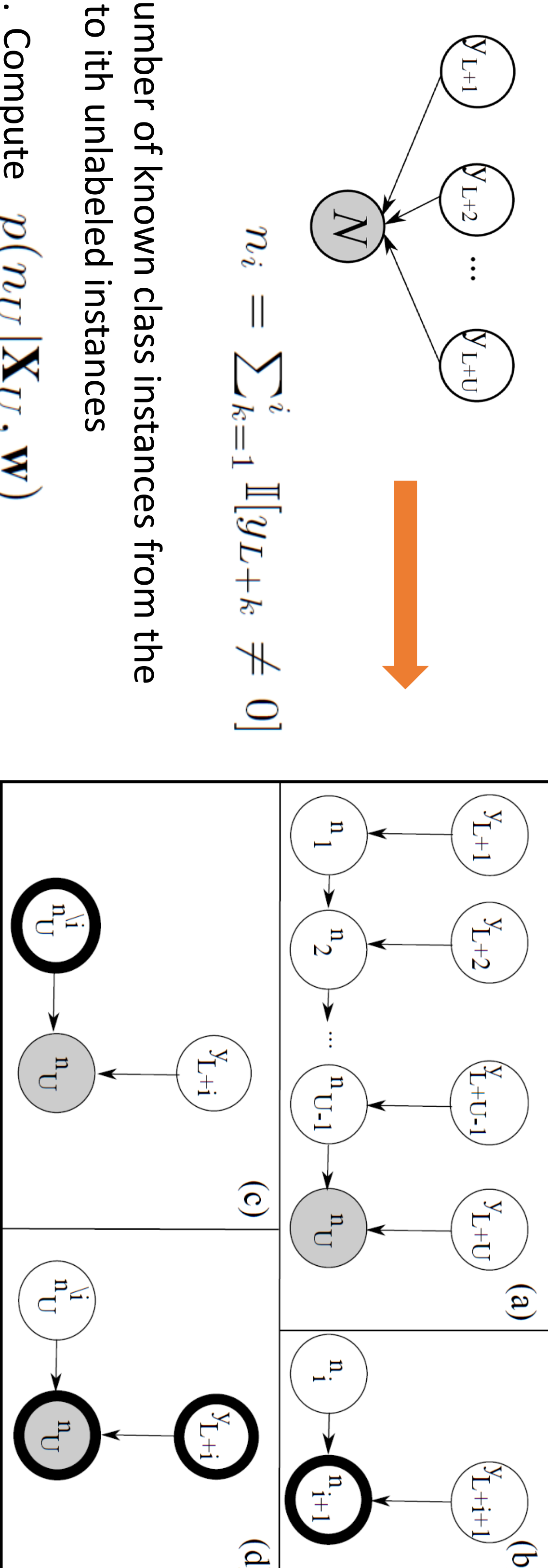
5 Inference

The log-likelihood $L(w) = \sum_{i=1}^L \log p(y_i | x_i, w) + \log p(N | \mathbf{X}_U, w)$

The surrogate function $g(w, w') = E_{y_U | D, w'} \log p(y_U, y_U, N | \mathbf{X}_L, \mathbf{X}_U, w) + E_{y_U | D, w'} \log p(y_U | \mathbf{X}_L, w) + \log p(y_U | \mathbf{X}_U, w) + E_{y_U | D, w'} \log p(N | y_U, \mathbf{X}_U, w)$

$$g(w, w') = \sum_{i=1}^L \sum_{c=0}^C \mathbb{I}[y_i = c] w_c^T x_i - \log \left(\sum_{c=0}^C e^{w_c^T x_i} \right) + \sum_{i=1}^U \sum_{c=0}^C p(y_{L+i} = c | N, \mathbf{X}_U, w') [w_c^T x_{L+i} - \log \left(\sum_{l=0}^C e^{w_l^T x_{L+i}} \right)]$$

6 E-step: Compute $p(y_{L+i} = c | N, \mathbf{X}_U, w)$



7 2D Toy dataset results

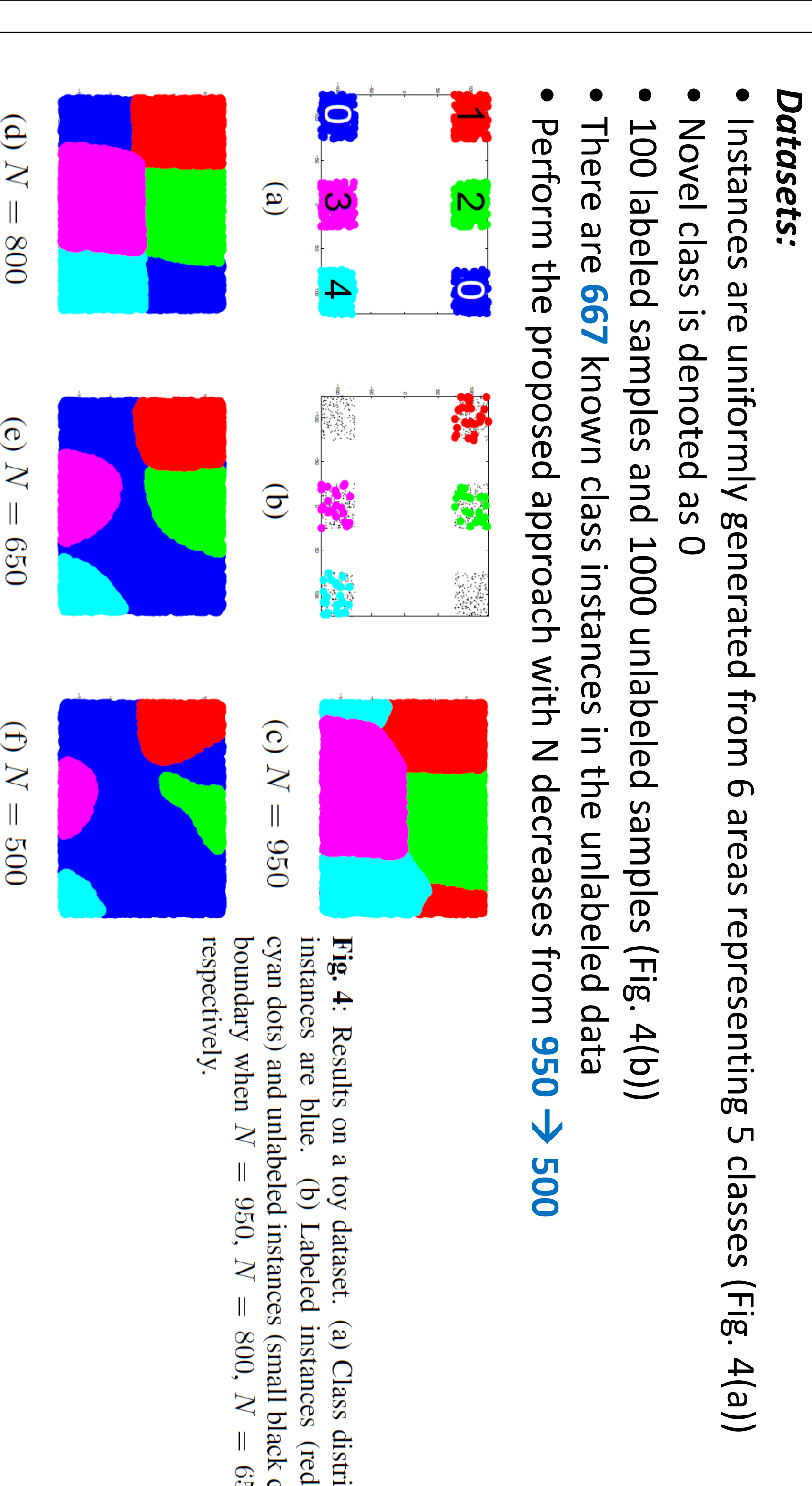


Fig. 4: Results on a toy dataset. (a) Class distribution, novel class instances are blue. (b) Labeled instances (red, pink, green, and cyan dots) and unlabeled instances (small black dots). (c-f) Learned boundary when $N = 950$, $N = 800$, $N = 650$, and $N = 500$, respectively.

8 Real dataset results

Baseline methods:

1. LRSSS-opt: proposed method
2. LRSSS-tune: proposed method with tuned N
3. LACU [1]
4. LR-L: trained with labeled data only
5. LRSSS-true: trained with labeled data only

Datasets: H1A bird [2], MSCV2, 50Salad, MNIST: Using the default class order: class1, 2, 3 are known, class4 is novel

Tuning N :

- Select labeled data as validation set
- Accuracy vs N curve
- Detect the knee

Dataset	H1A bird	MSCV2	50Salad	MNIST
LRSSS-opt	74.4±1.7	77.8±2.2	72.8±2.2	82.0±3.1
LRSSS-tune	73.2±2.5	69.0±1.5	69.0±2.7	69.2±3.3
LACU	51.6±8.8	65.8±7.1	65.7±7.5	79.4±3.3
LR-L	54.1±2.1	73.2±1.7	42.1±1.9	67.7±2.4
LRSSS-true	74.3±1.3	78.0±2.1	70.1±2.5	78.5±3.3

Table 1: Accuracy results of the proposed approach and LACU. The proposed method and indistinguishable values using 95% confidence two-tailed paired t -tests with the highest values are bolded.

9 Conclusion

- A framework for semi-supervised learning in the presence of novel class instances
- A probabilistic approach to control the cardinality of known class instances in the unlabeled data
- Future work:
 - Estimate the number of novel class instances
 - Improve the computational efficiency for the E-step

References

- [1] Q. Da, Y. Yu, and Z.-H. Zhou, "Learning with augmented class by exploiting unlabeled data," in AAAI Conference on Artificial Intelligence, 2014, pp. 1760–1766.