

# Pareto Optimal Streaming Unsupervised Classification

Sanjay Shakkottai<sup>1</sup>

Joint with Soumya Basu<sup>1</sup>,  
Steven Gutstein<sup>2</sup>, Brent  
Lance<sup>2</sup>

1. UT Austin, 2. ARL

# Streaming Unsupervised Learning

**Tasks:** Images arriving as a stream for classification

**Agents:** Human agents, Neural-net classifiers

- **Different expertise:** Unknown confusion matrices
- **Frozen labeling:** Label for each image fixed
- **Processing Rate:** One image/round

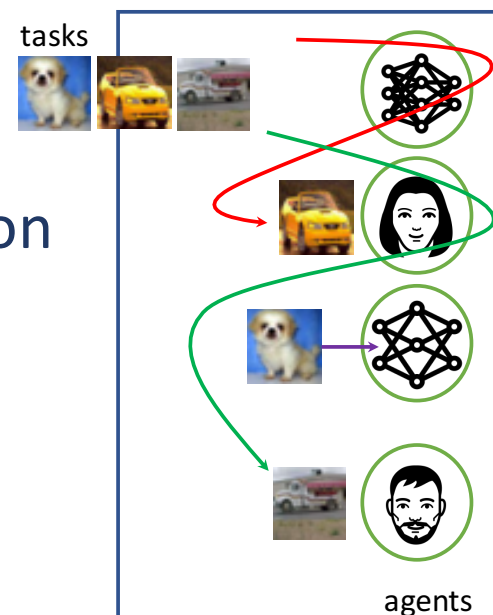


Image credits: CIFAR-10, A. Krizhevsky, 2009; thenounproject.com, (NNs - K. M. Synstad; Faces - A. Selimov)

## Types of Animals



Cat

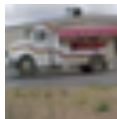


Dog

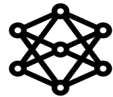
Car



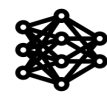
Truck



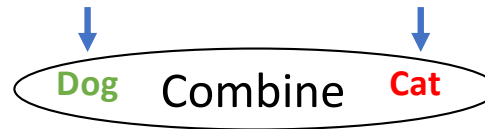
## Types of Vehicles



Animal  
Expert



Vehicle  
Expert



Dog

**Safest Strategy:**  
Send each image (task)  
to both the Classifiers

High accuracy, but low throughput

## Types of Animals



Cat

Dog

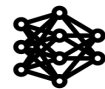
Car

Truck



Animal Expert

Car



Vehicle Expert

Dog

**Fastest Strategy:**  
Send each image (task)  
randomly to one  
Classifier

## Types of Vehicles

Lower accuracy, but high throughput

## Types of Animals



Cat

Dog

Car

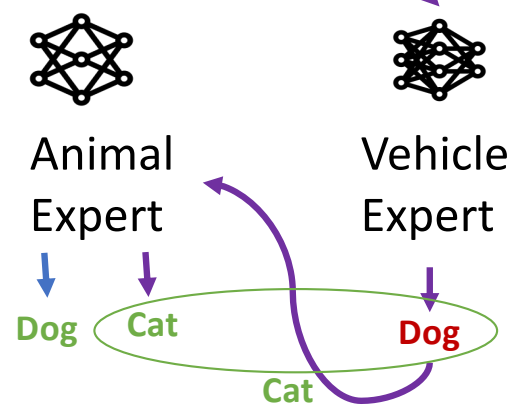
Truck



## Types of Vehicles

**Our Approach: Dynamic Routing**

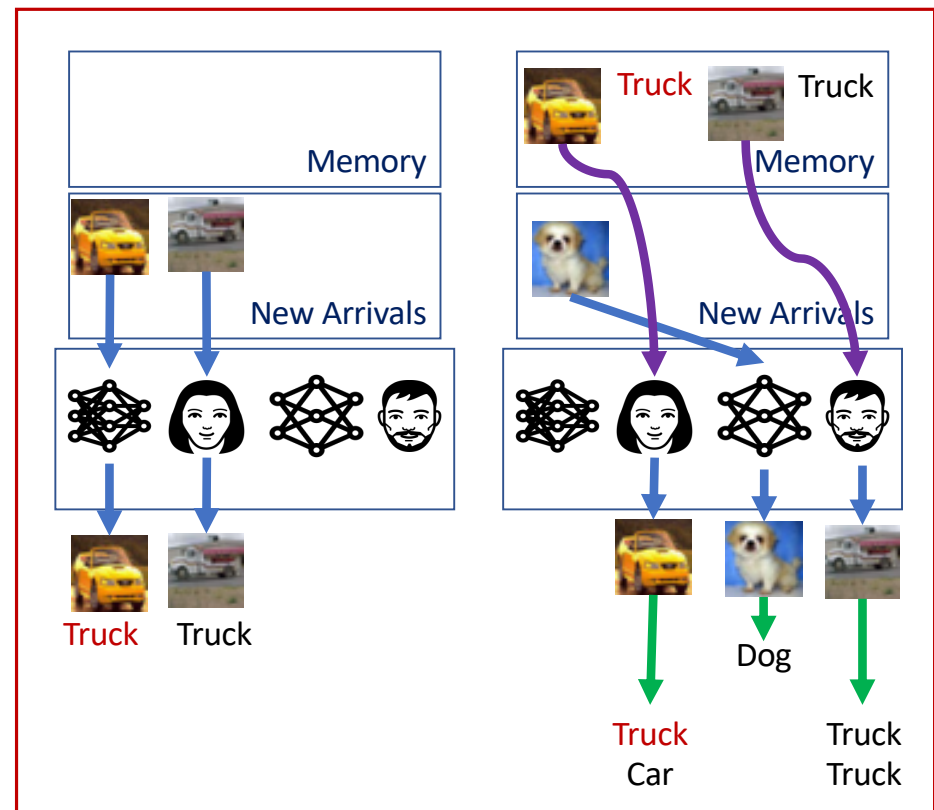
- Send randomly to one classifier
- If accuracy is low send to the remaining



Desired Accuracy  
High Speed

# Resource Allocation and Label Aggregation

- Sequential routing to subset of agents
- Collected labels are continually aggregated
- Aggregation using estimated confusion matrices



# Streaming Dawid-Skene Model

- Time slotted system,  $\mathbf{K}$  classes of images, and  $\mathbf{M}$  classifiers
- Independent identically distributed arrivals with rate  $\lambda$  **images/timeslot**
- Classifier characterized by its confusion matrix

A classifier's labels for a specific image is **frozen**. Repeatedly sending an image to the same classifier does not result in new labels.

Example: Trained Neural Network

# Pareto Optimality

- **Causal Policy:**  
Routing and aggregation are randomized functions of observed history
- **Confidence** =  $\mathbb{P}[\text{True label} = \text{final label} \mid \text{History}]$

## Threshold Accuracy ( $\theta$ )

Accurately labeled with Confidence  $\geq \theta$

OR

Sent to ALL the classifiers

Final label =  $f(\text{History, independent randomness})$

- **Arrival Rate ( $\lambda$ ) vs Threshold Accuracy ( $\theta$ ):**  
Higher  $\theta$  requires more classifiers per image on average  
Higher the threshold  $\theta$  the lower the arrival rate we can support
- **Goal:** Achieve the “best” trade-off between throughput ( $\lambda$ ) and accuracy ( $\theta$ )



# Prior Work: One-shot Unsupervised Learning

- **Batch of samples** without knowledge of true label (no ground truth)
- **Fixed classifiers** with unknown confusion matrices
- **Aim:** Combine label of **all** the classifiers
- **EM:** Dawid et al. JRSS'1979, Liu et al. Neurips'12, Zhang et al. Neurips'14
- **Majority Voting:** Li et al. Stat'14, Parisi et al. PNAS'14
- **Learning Confusion Matrices:** Zhang et al. Neurips'14, Jain et al. COLT'14
  
- **Differences - In our setting:**
  - Streaming classification
  - Rate vs accuracy tradeoff

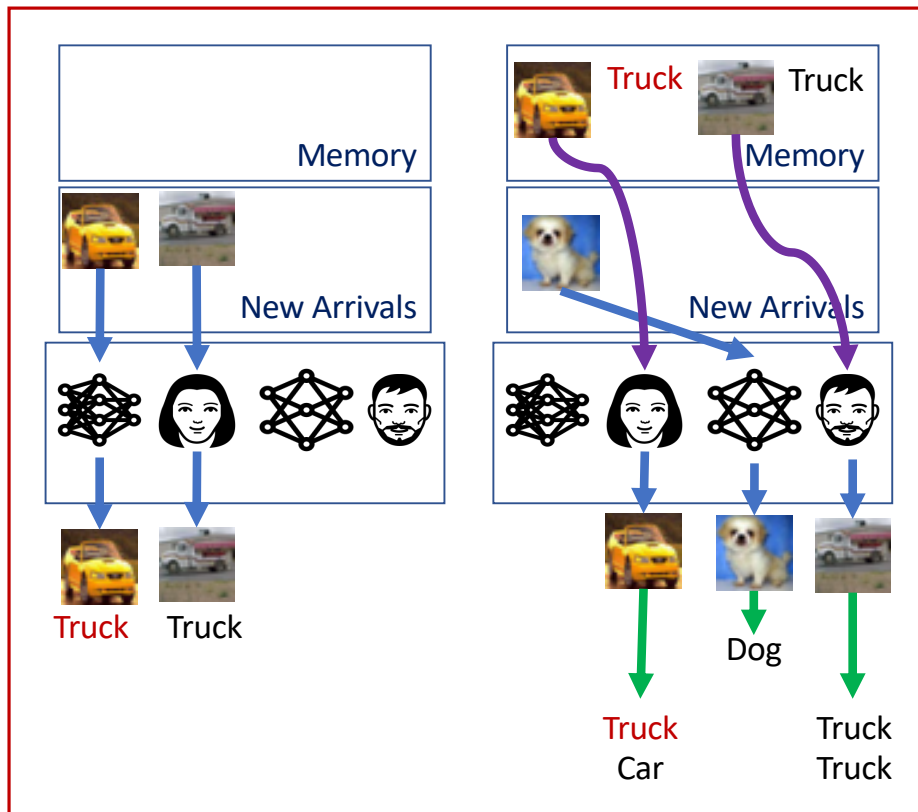
# Prior Work: Budgeted Crowdsourcing

- **Finite batch** of samples at beginning of time
- Classifiers from crowd arrive online (**effectively randomized classifiers**)
- Fixed budget of total number of classifier actions (over all samples)
- **Aim:** Achieve maximum accuracy given budget
- **Optimal accuracy rates in budgeted crowdsourcing:**
  - Karger et al. Neurips'11, Khetan et al. Neurips'16
- **Differences - In our setting:**
  - Samples arrives online
  - Frozen classifiers

# Prior Work: Information Processing Networks

- Fixed but **randomized** processors
- Accuracy can be arbitrarily high
- Streaming arrival of samples
- **Exogenous** departure from processors
- **Aim:** Support maximum arrival rate
- **Capacity region characterization:** Shah et al. Allerton'17
  
- **Differences - In our setting:**
  - Endogenous departures
  - Frozen classifiers

# How Does Belief Evolve? (1/2)



$$\overbrace{Events(t)}^{Schedule(t) \rightarrow Labels(t)}$$

$$History(t) = \bigcap_{s=1}^t Events(s)$$

$$\mathbb{P}[\bigcap_{image\ j} T_j = \mathbf{k}(j) | History(t)]$$

Belief at time t

$T_j$ : true label of image j

## How Does Belief Evolve? (2/2)

$\mathbf{T}_j$ : true label of image  $j$

$L(i, j)$ : Label of image  $j$  from classifier  $i$

$Cl(j, t)$ : Classifiers assigned to image  $j$  up to time  $t$

Is this true?

$$\mathbb{P}[\cap_{image\ j} \mathbf{T}_j = \mathbf{k}(j) | \mathbf{History}(t)] = \frac{1}{Z'} \prod_{image\ j} \mathbb{P}[\mathbf{T}_j = \mathbf{k}(j) | \cap_{i \in Cl(j,t)} L(i, j) = l(i, j)]$$

The conditional belief on the true label an image potentially depends on entire history of all co-existing images and their labels

- The **set of labels** acquired by an image until time  $t$  depends on the **set of classifiers** that the image was scheduled to, until time  $t$
- This **set of classifiers** depends on the **scheduling decisions** until time  $t$
- These **scheduling decisions** depends on the entire past history of **labels of all images** in the system until time  $t$

# Product Form of Belief Evolution

- Distribution of history depends on the the true label only through the collected labels of the sample

$\mathbf{T}_j$ : true label of image j

$L(i, j)$ : Label of image j from classifier i

$Cl(j, t)$ : Classifiers assigned to image j up to time t

## Product form on the Belief Evolution

$$\begin{aligned}\mathbb{P}[\cap_{image\ j} \mathbf{T}_j = \mathbf{k}(j) | \mathbf{History}(t)] &= \frac{1}{Z'} \prod_{image\ j} \mathbb{P}[\mathbf{T}_j = \mathbf{k}(j) | \cap_{i \in Cl(j,t)} L(i, j) = l(i, j)] \\ &= \frac{1}{Z(t)} \prod_{image\ j} \prod_{i \in Cl(j,t)} \mathbb{P}[L(i, j) = l(i, j) | \mathbf{T}_j = \mathbf{k}(j)] \mathbb{P}[\mathbf{T}_j = \mathbf{k}(j)] \quad (\text{Bayes Form})\end{aligned}$$

# Product Form of Belief Evolution

$$\begin{aligned} \text{Induction on } t: & \mathbb{P}[\mathbf{History}(t) | \cap_j T_j = k(j)] \\ &= \frac{1}{Z'(t)} \prod_j \prod_{i \in \text{cl}(j,t)} \mathbb{P}[L(i,j) = l(i,j) | T_j = k(j)] \end{aligned}$$

$$\overbrace{\text{Sch}(t) \rightarrow \text{Labels}(t)}^{\text{Events}(t)}$$

$$\begin{aligned} \mathbb{P}[\mathbf{History}(t+1) | \cap_j T_j = k(j)] &= \mathbb{P}[\mathbf{Events}(t+1) \cap \mathbf{History}(t) | \cap_j T_j = k(j)] \\ &= \mathbb{P}[\mathbf{Events}(t+1) | \mathbf{History}(t), \cap_j T_j = k(j)] \mathbb{P}[\mathbf{History}(t) | \cap_j T_j = k(j)] \\ &= \mathbb{P}[\text{Sch}(t+1) | \mathbf{History}(t), \cap_j T_j = k(j)] \times \mathbb{P}[\text{Labels}(t+1) | \cap_j T_j = k, \text{Sch}(t)] \\ &\quad \times \frac{1}{Z'(t)} \prod_j \prod_{i \in \text{cl}(j,t)} \mathbb{P}[L(i,j) = l(i,j) | T_j = k(j)] \end{aligned}$$

Causal Policy

continued...

# Product Form of Belief Evolution

$$\begin{aligned} & \text{Induction on } t: \mathbb{P}[\mathbf{History}(t) | \cap_j T_j = \mathbf{k}(j)] \\ &= \frac{1}{Z'(t)} \prod_j \prod_{i \in \mathbf{cl}(j,t)} \mathbb{P}[L(i,j) = l(i,j) | T_j = \mathbf{k}(j)] \end{aligned}$$

$$= \mathbb{P}[\cap_{(i,j) \in \mathbf{Sch}(t+1)} L(i,j) = l(i,j) | \cap_j T_j = \mathbf{k}(j)]$$

*Sch*(*t*) is a set of (classifier, image) pairs

$$\times \frac{\mathbb{P}[\mathbf{Sch}(t+1) | \mathbf{History}(t)]}{Z'(t)} \prod_j \prod_{i \in \mathbf{cl}(j,t)} \mathbb{P}[L(i,j) = l(i,j) | T_j = \mathbf{k}(j)]$$

$$= \prod_{(i,j) \in \mathbf{Sch}(t+1)} \mathbb{P}[L(i,j) = l(i,j) | T_j = \mathbf{k}(j)]$$

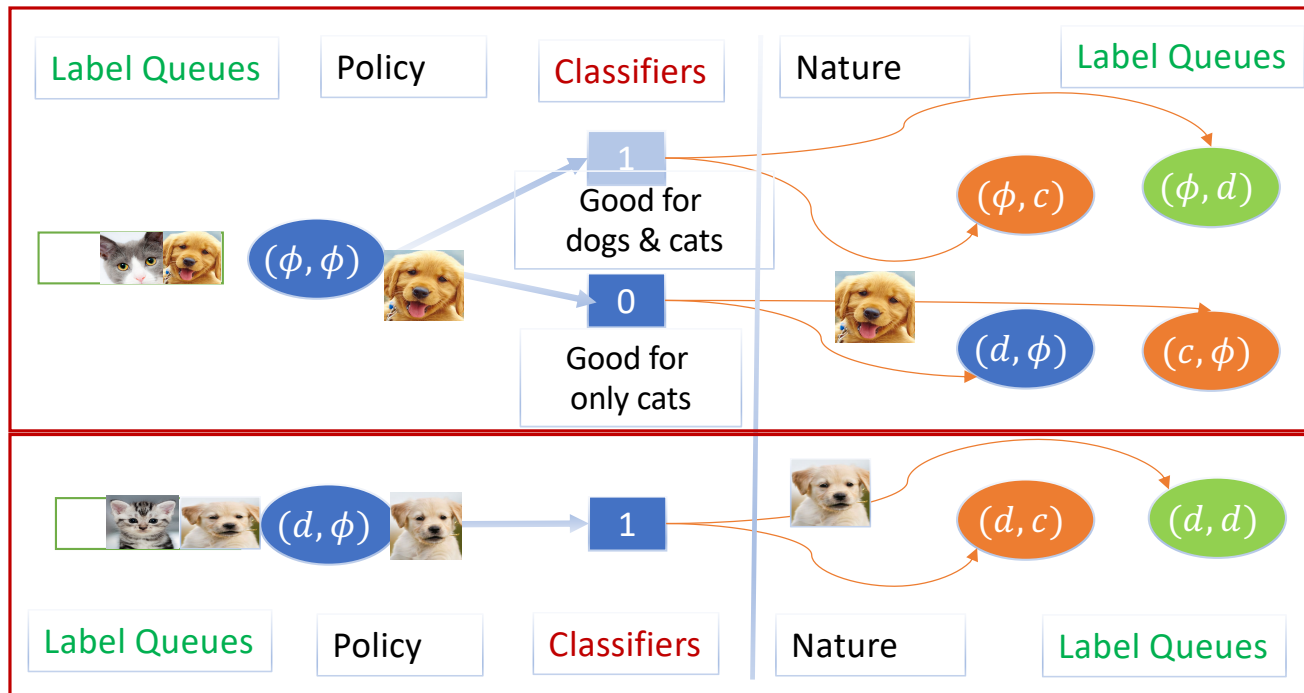
After scheduling, new label of image *j* depends only on true label of *j*

$$\times \frac{\mathbb{P}[\mathbf{Sch}(t+1) | \mathbf{History}(t)]}{Z'(t)} \prod_j \prod_{i \in \mathbf{cl}(j,t)} \mathbb{P}[L(i,j) = l(i,j) | T_j = \mathbf{k}(j)]$$



# Network of Label Queues

- **Partial Label:** Observed label tuple, e.g. (cat, none), (dog, cat), (none, none)
- **Label Queues:** Binary classification (dogs vs cats) with 2 classifiers



Example 1

Example 2

Policy Controls **blue edge**  
Nature controls **red edge**

**Scheduling:**  
Classifiers and Label Assignment

**Departure Labels:**  
Labels for which we have desired accuracy or all classifiers are used



# Expected MaxWeight for Agent Selection

- Let  $Q_\ell(t)$  be queue length of label queue  $\ell$  at time  $t$
- Under assignment  $\mathbf{S}$ 
  - Departure from queue  $\ell$ ,  $D_\ell(\mathbf{S})$
  - Arrival into queue  $\ell$ ,  $A_\ell(\mathbf{S})$
- Choose assignment according to Expected MaxWeight:

$$\mathbf{S}(t) = \underset{\mathbf{S}}{\operatorname{argmax}} \sum_{\ell} Q_\ell(t) (D_\ell(\mathbf{S}) - \mathbb{E}[A_\ell(\mathbf{S})])$$

- Belief is only dependent on the collected labels
- Expectation on Nature's choice does not vary over time
  - Depends only on labels and confusion matrix

# Performance Guarantees

- **Pareto region:**

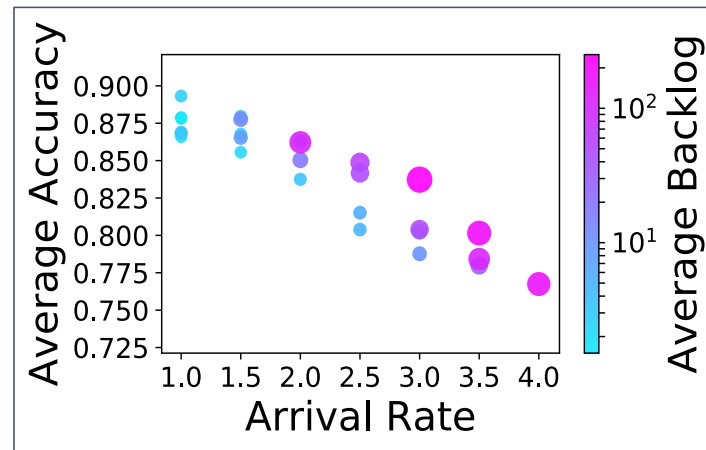
Set of ( Arrival rate ( $\lambda$ ) , threshold ( $\theta$ ) ) tuple so that the sum of expected queue length remains bounded under some causal policy

- Pareto region is characterized by a **maximum network flow** problem

- **Pareto Optimality:**

For any arrival rate  $\lambda$  and threshold  $\theta$ , for any  $\epsilon > 0$ ,  
if  $(\lambda + \epsilon, \theta + \epsilon)$  lies in the Pareto region  
then **Expected Backpressure + online tensor decomposition**  
support  $(\lambda, \theta)$  with bounded expected memory.

# Neural Network Ensembles

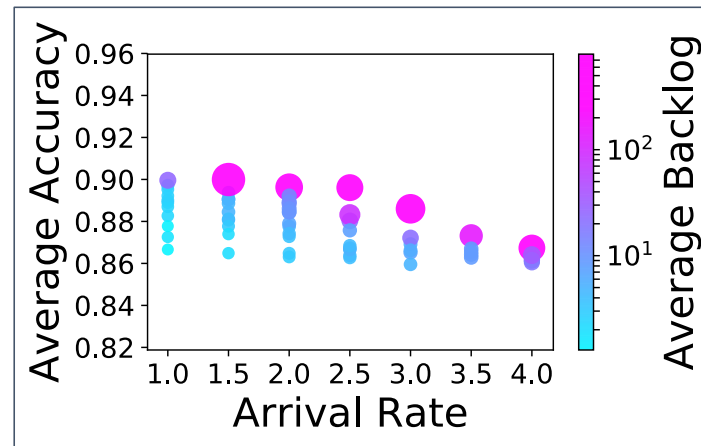


## Experiment 1

### Pareto Region Exp 1

- **6 Classifiers:** Three AlexNet, one VGG-19, and two ResNet-18
- **3 Labels:** Group 1 (airplanes, ships, trucks, cars), Group 2 (birds, frogs cats), and Group 3 (dogs, deer, horses)

# Neural Network Ensembles



**Improved Classifiers  
lead to improved  
Pareto Region**

## Experiment 2

### Pareto Region Exp 2

- **6 Classifiers:** One VGG-11, one VGG-16, two VGG-19, and two ResNet-18
- **5 Labels:** Group 1 (airplanes, ships), Group 2 (trucks, cars), Group 3 (birds, frogs), Group 4 (cats, dogs) and Group 5 (deer, horses)

# Future Directions

- What fraction of Pareto Region can be covered by simple algorithms?
  - Algorithms without explicitly recovering parameters  
e.g. Majority voting, Routing over limited number of rounds
- Can we maximize threshold accuracy given an arrival rate?
- Can we maximize average accuracy inside Pareto Region?

Backup Slides

# Online Learning of Confusion Matrices

- Explore (w.p.  $1/t$  at time  $t$ ) by sending one sample to all classifiers
- An  $(\alpha, \beta)$ -oracle, with  $n$  exploration samples, outputs confusion matrices and true probability vector with  $L_\infty$  error at most  $O(n^{-\alpha})$  with probability at least  $1 - \Omega(n^{-\beta})$
- $(0.5 - \epsilon, 1 - \epsilon)$ -oracle created by adapting tensor decomposition based one-shot unsupervised learning [Zhang et al. 2014]
  - Reinitialize only when a new explore sample is obtained
  - Ensure w.p. 1 the initialization is 'good' for tensor power method
  - Iteratively improve estimates between two exploration instances



# Online Dawid-Skene Model

- Time slotted system,  $\mathbf{K}$  classes of images, and  $\mathbf{M}$  classifiers
- I.i.d. Arrival with rate  $\lambda$  images/timeslot
- **Classifier  $i$  :**
  - Confusion matrix  $\mathbf{C}_i(\cdot, \cdot)$
  - **Speed** = 1 image / timeslot
- **Image  $j$  :**  $(M+1)$ -tuple  $(L_j(1), \dots, L_j(M), T_j)$ 
  - **True label:** Label  $T_j \sim p_g \in \Delta_{\mathbf{K}}$  chosen and fixed
  - **Label from Classifier  $i$  :** Label  $L_j(i) \sim C_i(T_j, \cdot)$  chosen and fixed

A specific classifier labels a specific image deterministically  
e.g. Trained Neural Network

# Sufficiency of Compressed Causal Policy

- **Compressed History:**

Aggregate samples with identical collected labels, delete the sample ids

- **Compressed Causal Policy:**

Decisions are randomized functions of the compressed history

For any causal policy there is a compressed causal policy s.t.

- For all time  $t$ , the compressed history distribution is identical
- Threshold accuracy is attained is identical

**Key Proof Idea:**

The belief of a sample only depends on its collected labels.

The new labels and threshold accuracy only depend on the belief.