

Yandex

Crowdsourcing Natural Language Data at Scale: A Hands-On Tutorial

Alexey Drutsa, Dmitry Ustalov, Valentina Fedorova, Olga Megorskaya, Daria Baidakova

Part III:

Advanced Techniques

Dmitry Ustalov,
Analyst/Software Developer
Crowdsourcing Research Group

Tutorial Schedule

```
graph TD; Intro[Introduction: 15 min] --> PartI[Part I: 30 min  
Key Components for  
Data Collection]; PartI --> PartII[Part II: 60 min  
Practice Session I]; PartI --> PartIII[Part III: 45 min  
Advanced  
Techniques]; PartII --> Lunch[Lunch Break:  
45 min]; Lunch --> PartIV[Part IV: 30 min  
Practice Session II]; PartIV --> PartV[Part V: 15 min  
Conclusion];
```

Introduction: 15 min

Part I: 30 min
Key Components for
Data Collection

Part II: 60 min
Practice Session I

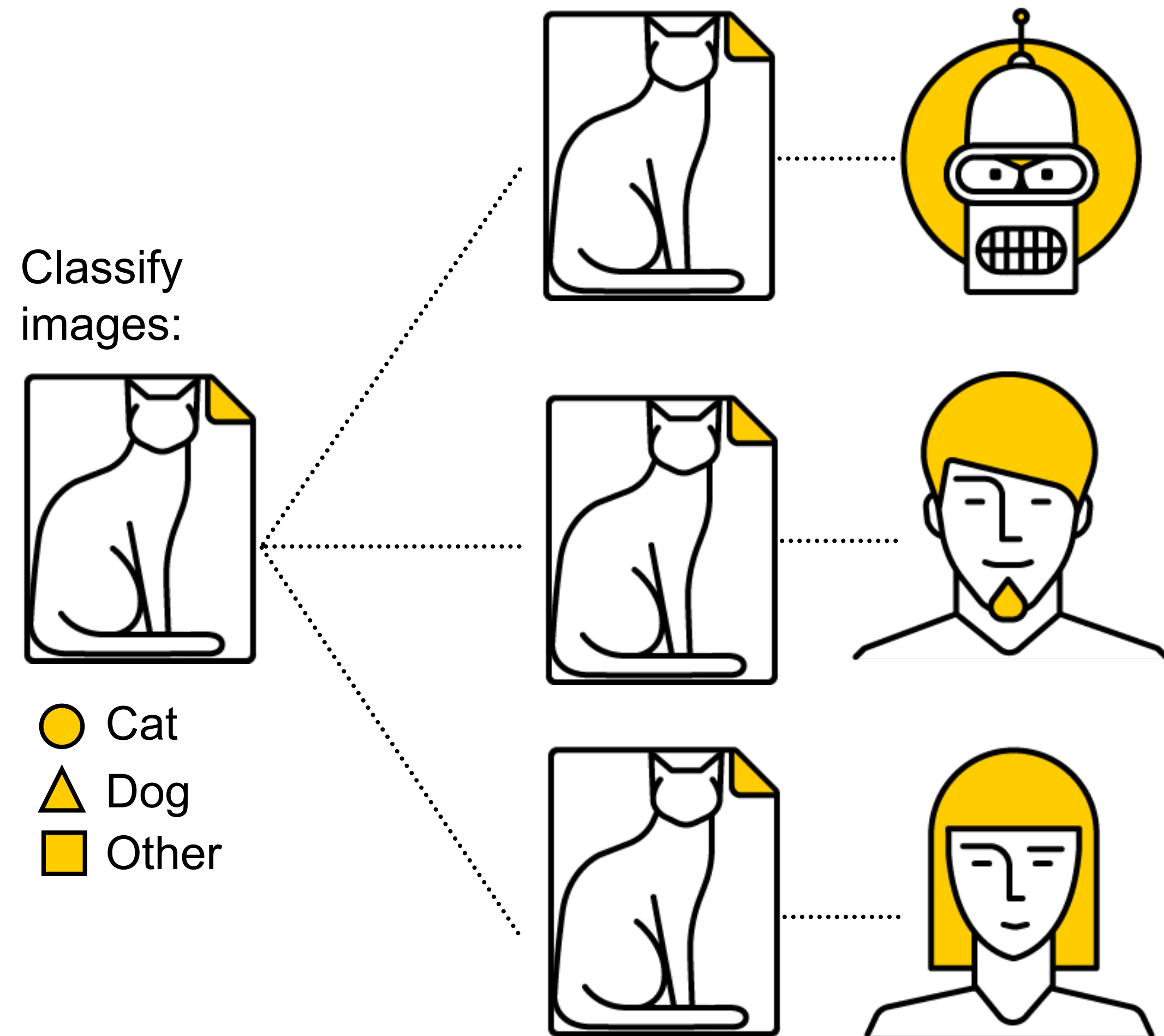
Lunch Break:
45 min

Part III: 45 min
Advanced
Techniques

Part IV: 30 min
Practice Session II

Part V: 15 min
Conclusion

Labeling data with crowdsourcing



› How to choose a reliable label?

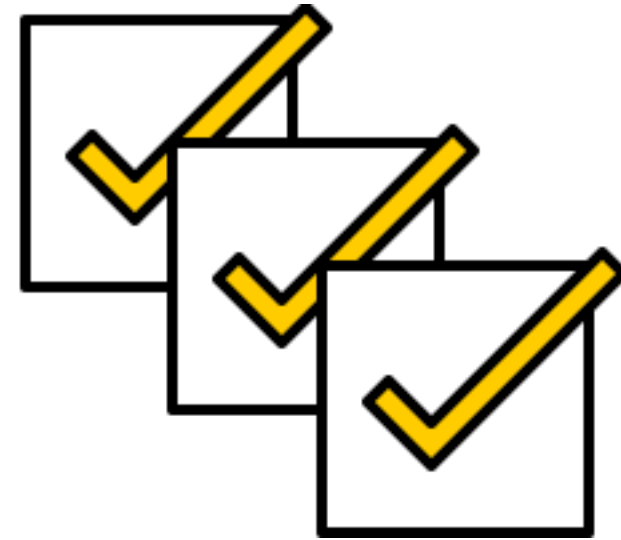
› How many labels per object?

› How much to pay for labels?

› ...

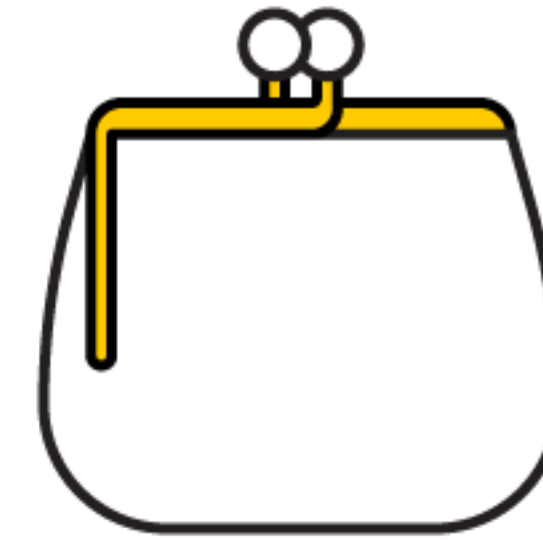
Evaluation of labeling approaches

Accuracy



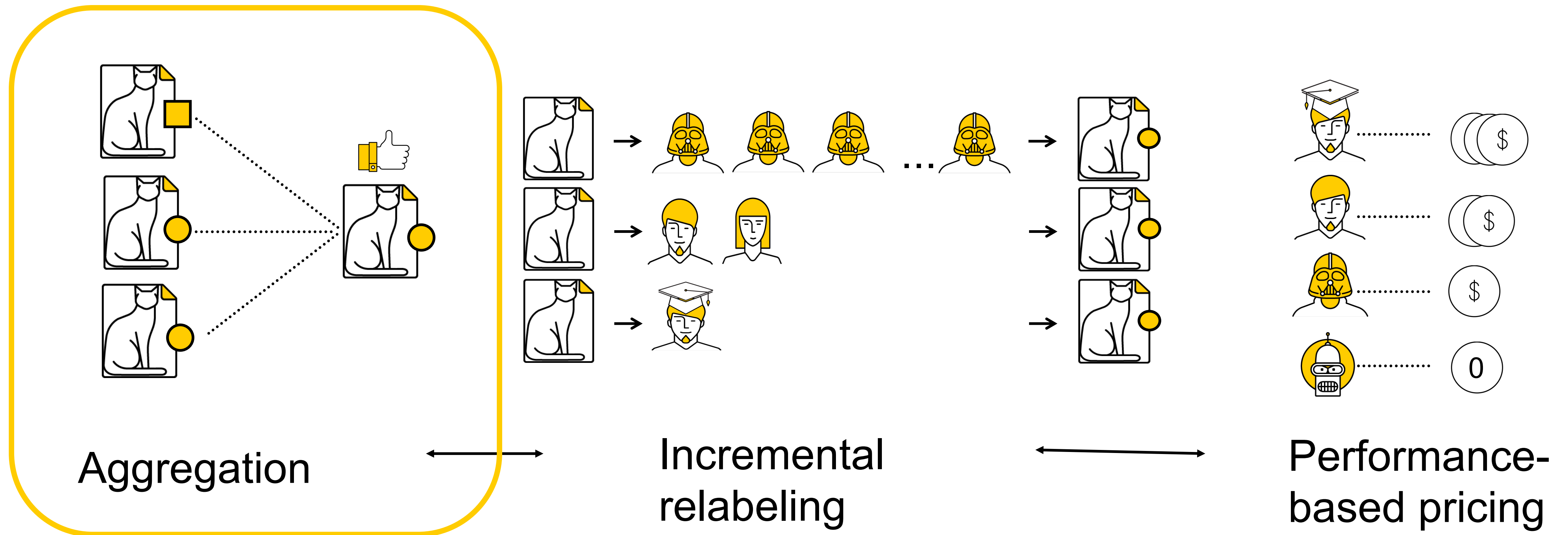
vs

Cost



- Labels with a maximal level of accuracy for a given budget
- or
- Labels of a chosen accuracy level for a minimal budget

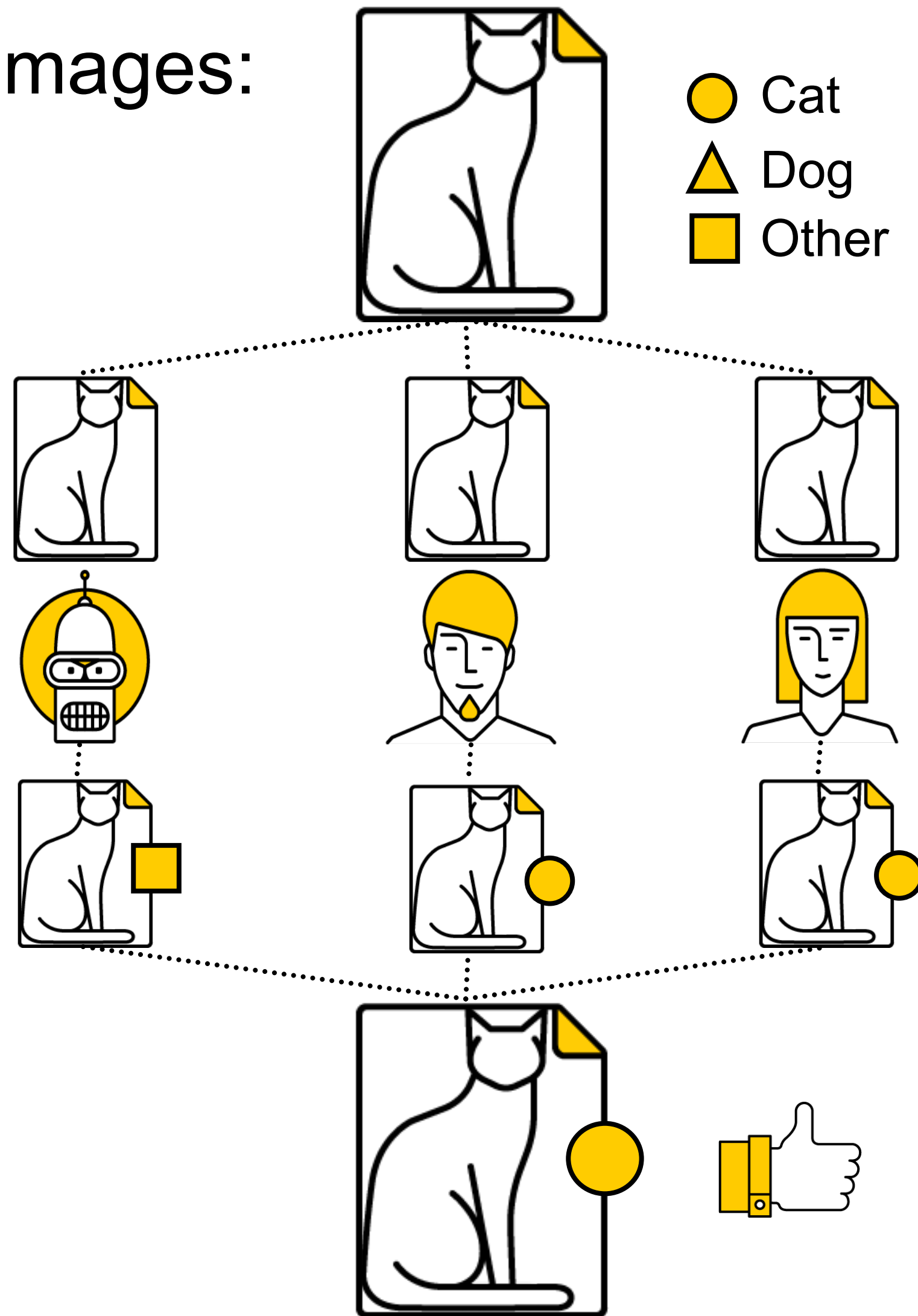
Key components of labeling with crowds



Aggregation

Labeling data with crowds

Classify images:



Upload multiple copies of each object to label

Performers assign noisy labels to objects

Aggregate multiple labels for each object into a more reliable one

Process results

Yandex Toloka

Projects

Users

Skills

Profile

Messages

?

\$0.00 / \$5.59

🇬🇧

Toloka requester

Projects > Does the image contains traffic lights? > pool

pool — closed

Statistics

Download results

View operations

Dawid-Skene aggregation model

Aggregation by skill

Edit

?

POOL TASKS (File example for task uploading (tsv, UTF-8)) ?

Upload

files

Edit

Preview

30 task suites

0 training task

90 tasks

10 control task

100 %

Done 30, accepted 30

View assignments

0

30

Multiclass labels

Multiclassification

Query: Machine learning
URL: https://en.wikipedia.org/wiki/Machine_learning

Open the original

Yandex

Google

1

☐

Vital

2

☐

Useful

3

☐

Relevant+

4

☐

Relevant-

5

☐

Irrelevant

6

☐

Not displayed



WIKIPEDIA
The Free Encyclopedia

Main page

Contents

Featured content

Current events

Random article

Donate to Wikipedia

Wikipedia store

Interaction

Help

About Wikipedia

Community portal

en.wikipedia.org Machine learning - Wikipedia

Not logged in Talk Contributions Create account Log in

Article Talk Read Edit View history Search Wikipedia

Machine learning

From Wikipedia, the free encyclopedia

For the journal, see [Machine Learning \(journal\)](#).
"Statistical learning" redirects here. For statistical learning in linguistics, see [statistical learning in language acquisition](#).

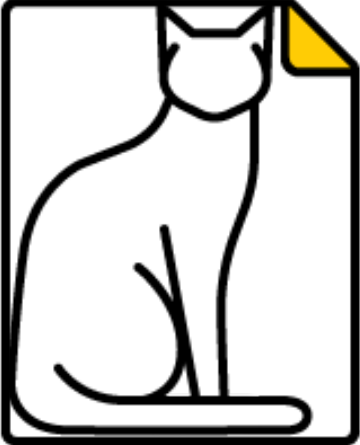
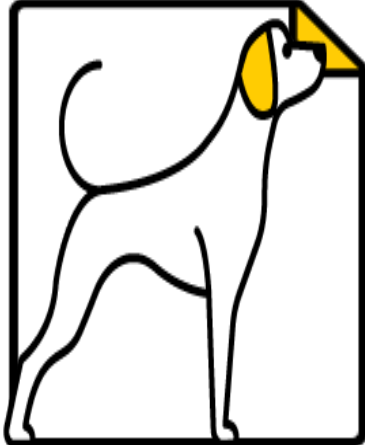

Machine learning (ML) is the scientific study of algorithms and statistical models that computer systems use in order to perform a specific task effectively without using explicit instructions, relying on patterns and inference instead. It is seen as a subset of artificial intelligence. Machine learning algorithms build a mathematical model based on

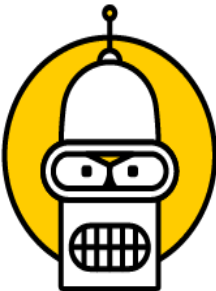
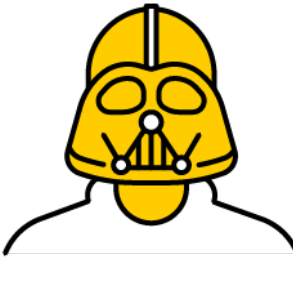



Machine learning and data mining



Notation

› Categories $k \in \{1, \dots, K\}$. E.g.: ● Cat ▲ Dog ■ Other

› Objects $j \in \{1, \dots, J\}$. E.g.:   

› Performers: $w \in \{1, \dots, W\}$. E.g.:     

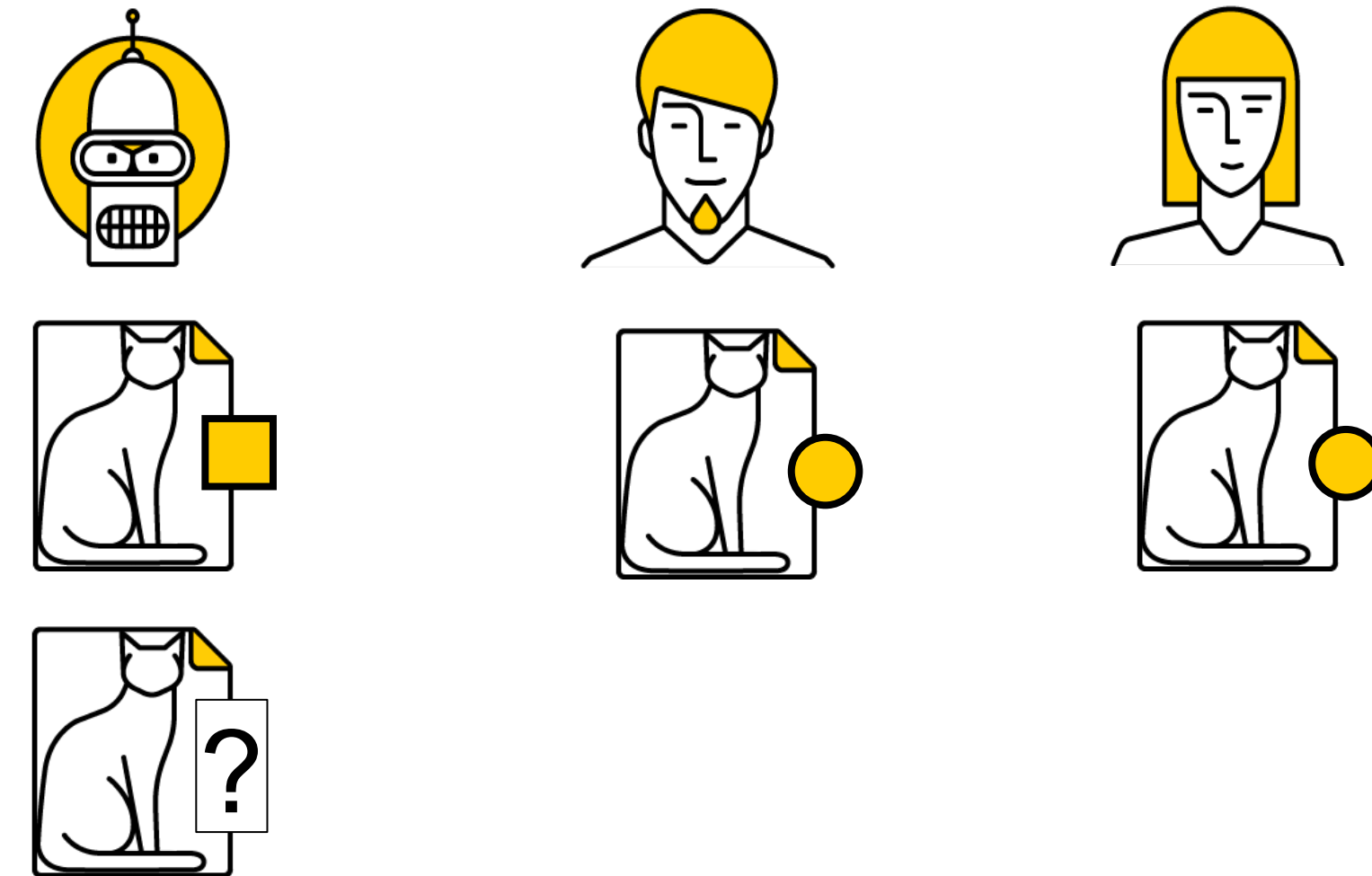
- $W_j \subseteq \{1, \dots, W\}$ - performers labeled object j

The simplest aggregation: Majority Vote (MV)

The problem of aggregation:

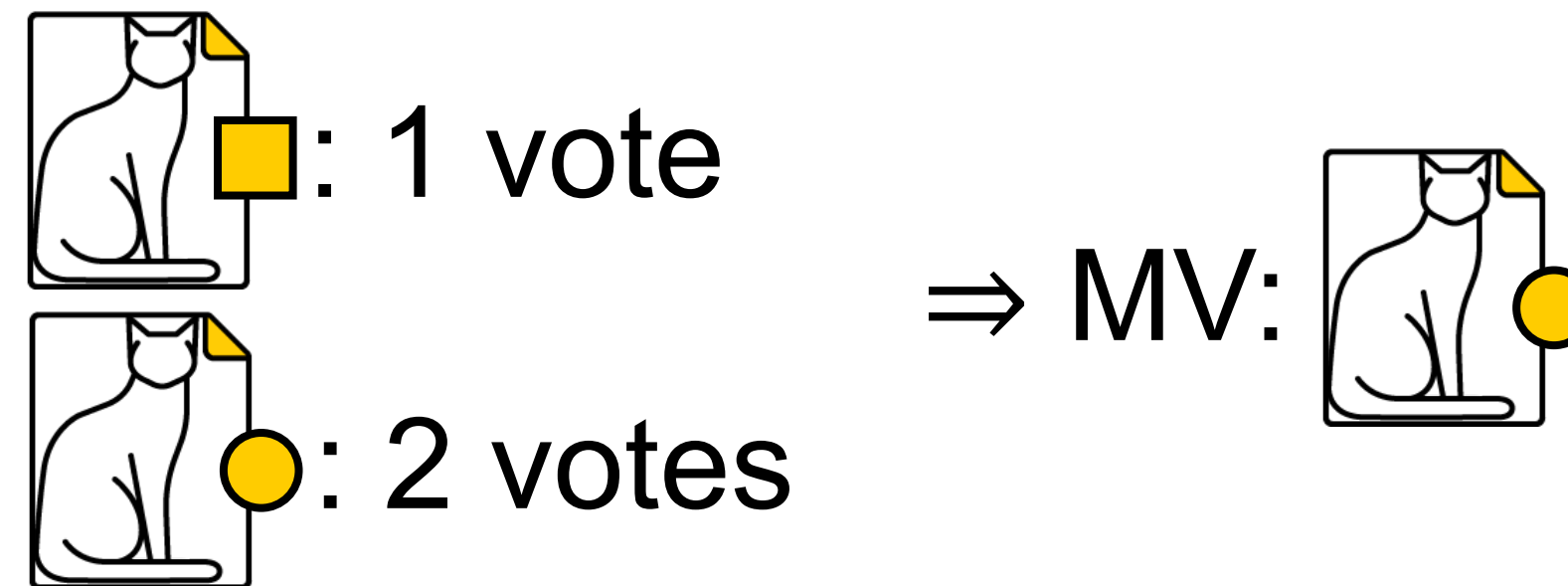
- › Observe noisy labels

$$\mathbf{y} = \{y_j^w \mid j = 1, \dots, J \text{ and } w = 1, \dots, W\}$$



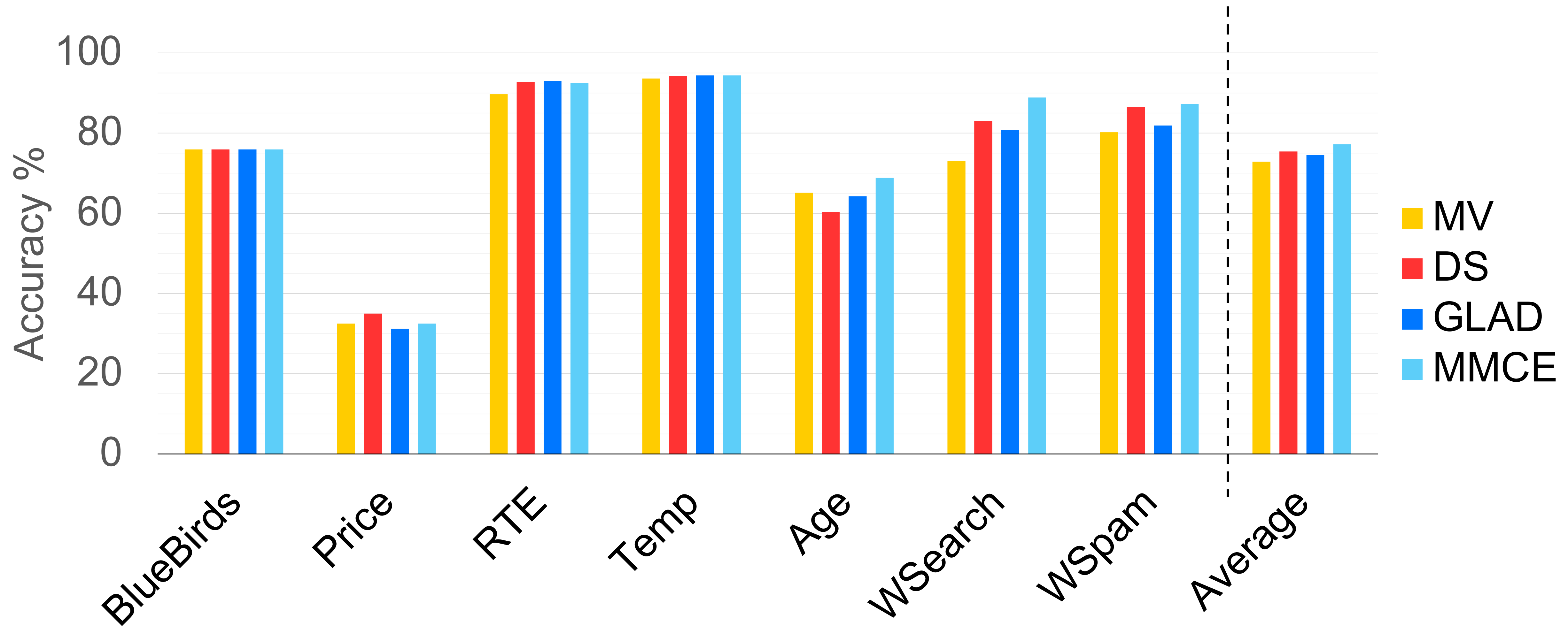
- › Recover true labels $\mathbf{z} = \{z_j \mid j = 1, \dots, J\}$

A straightforward solution:



$$\hat{z}_j^{MV} = \arg \max_{y=1, \dots, K} \sum_{w \in W_j} \delta(y = y_j^w), \text{ where } \delta(A) = 1 \text{ if } A \text{ is true and } 0 \text{ otherwise}$$

Performance of MV vs other methods

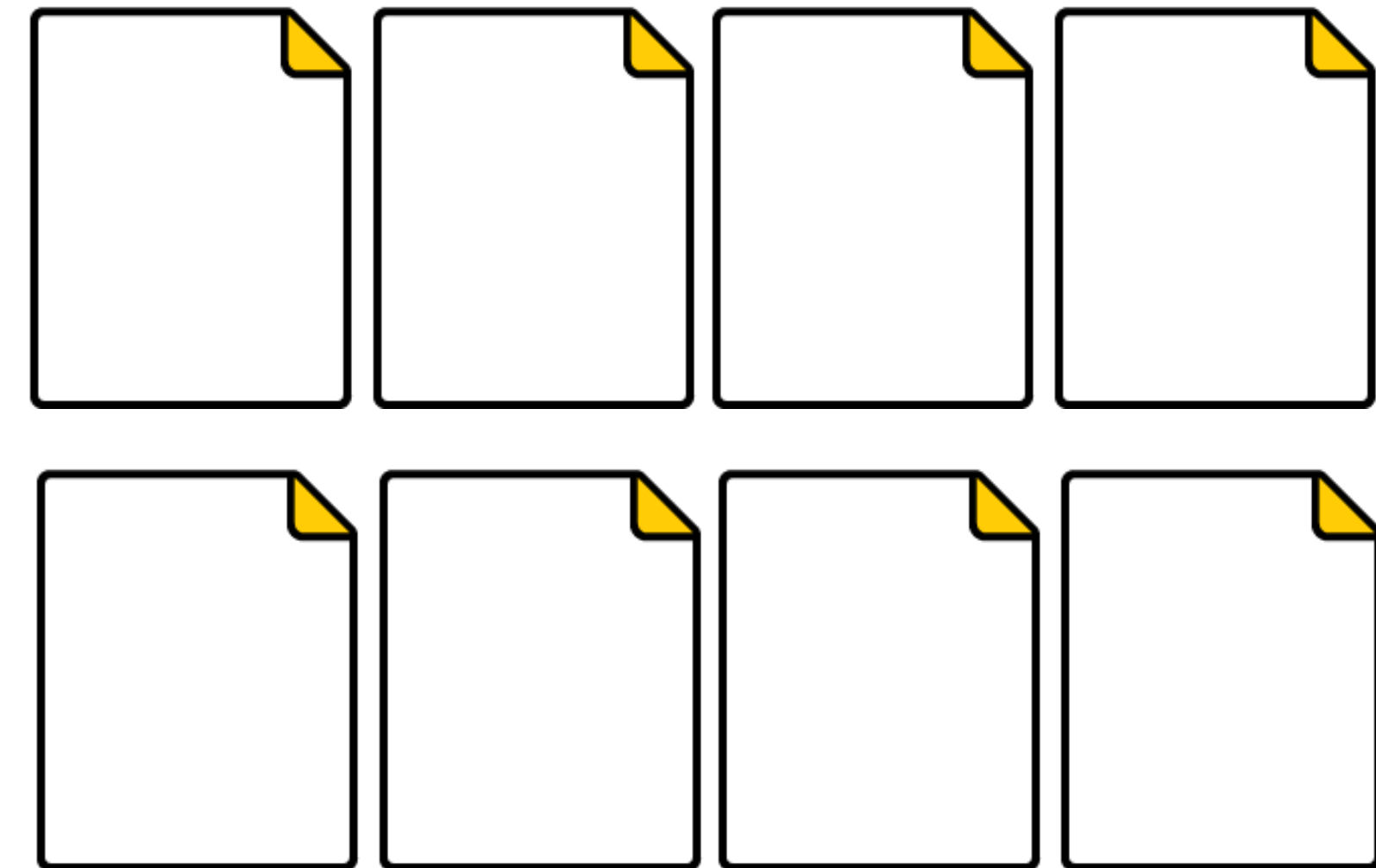


Properties of MV

› All performers are treated similarly



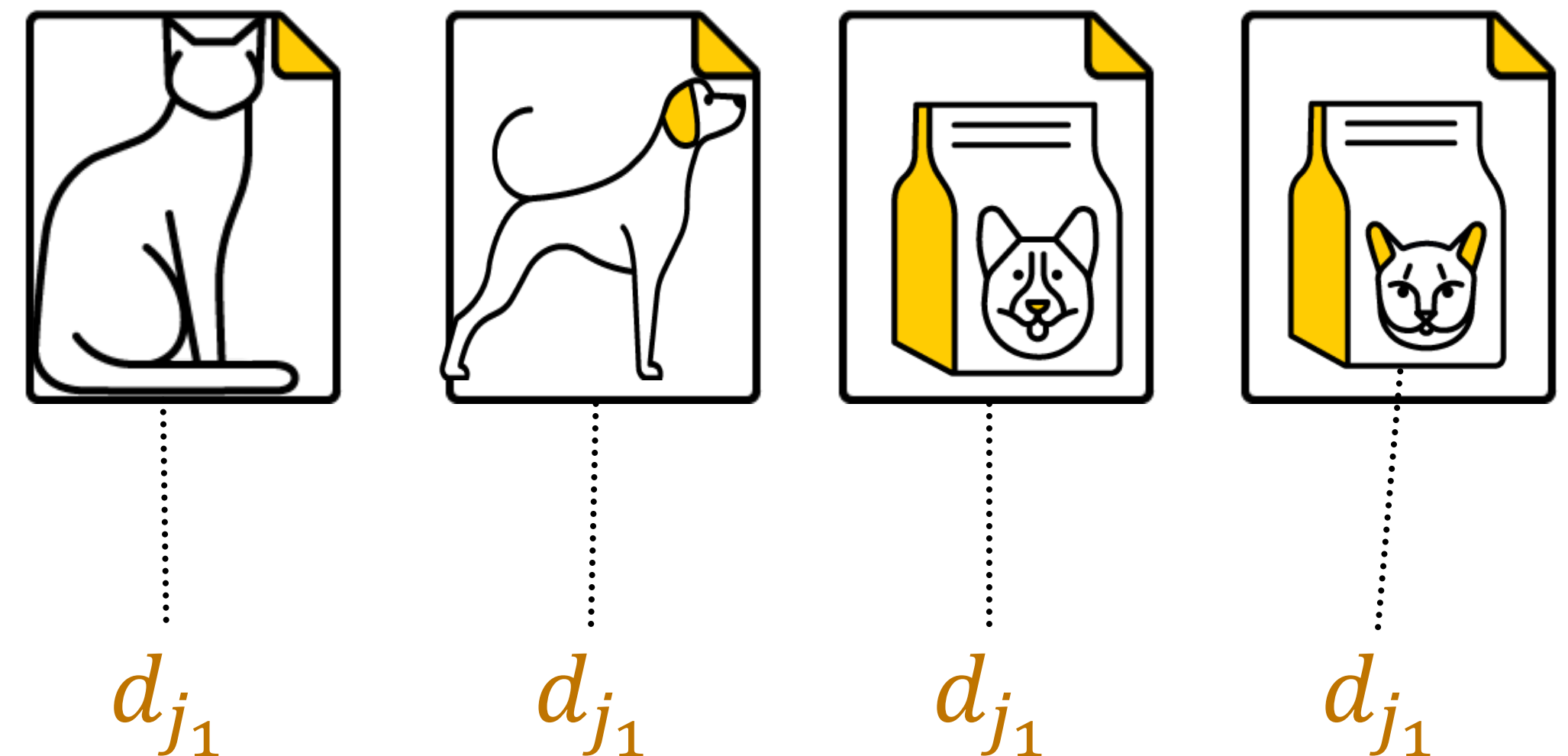
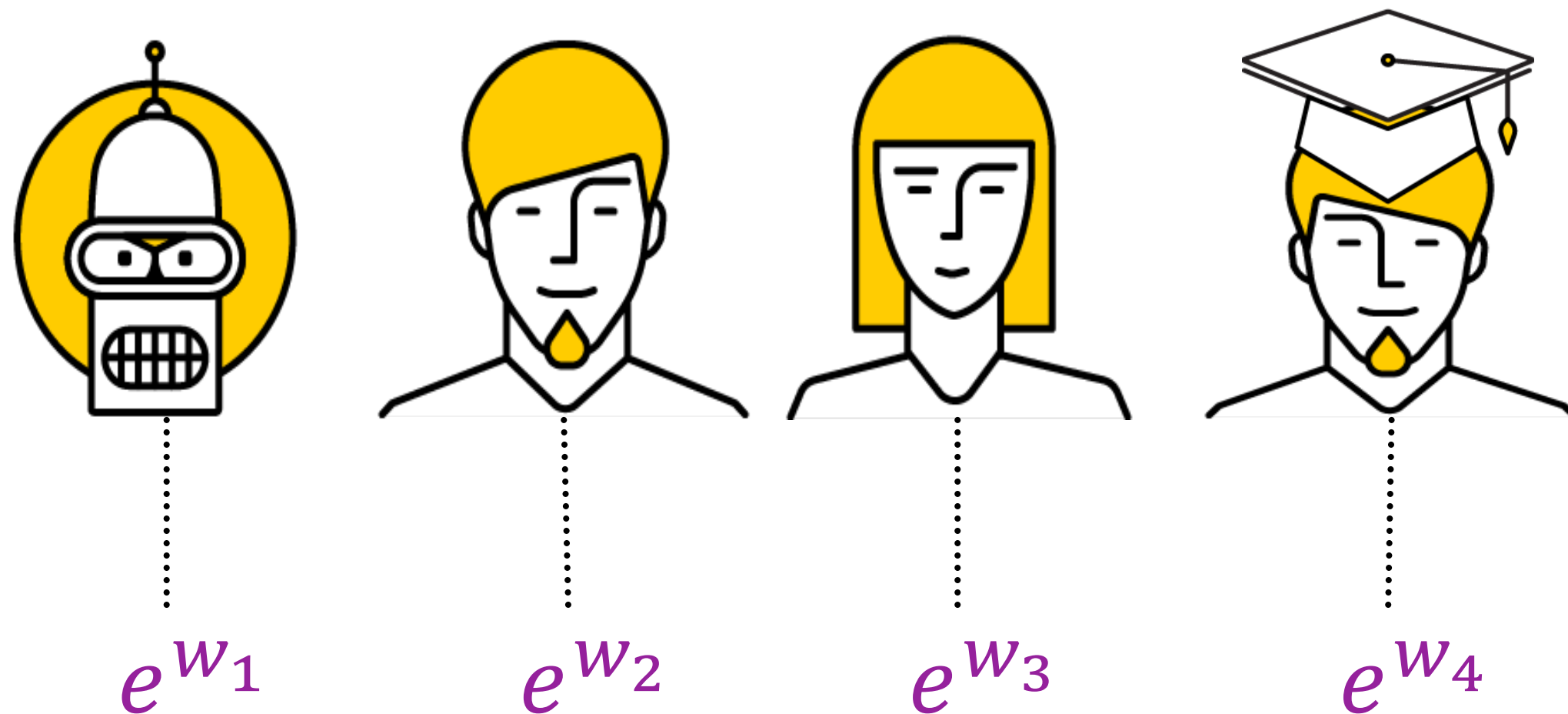
› All objects are treated similarly



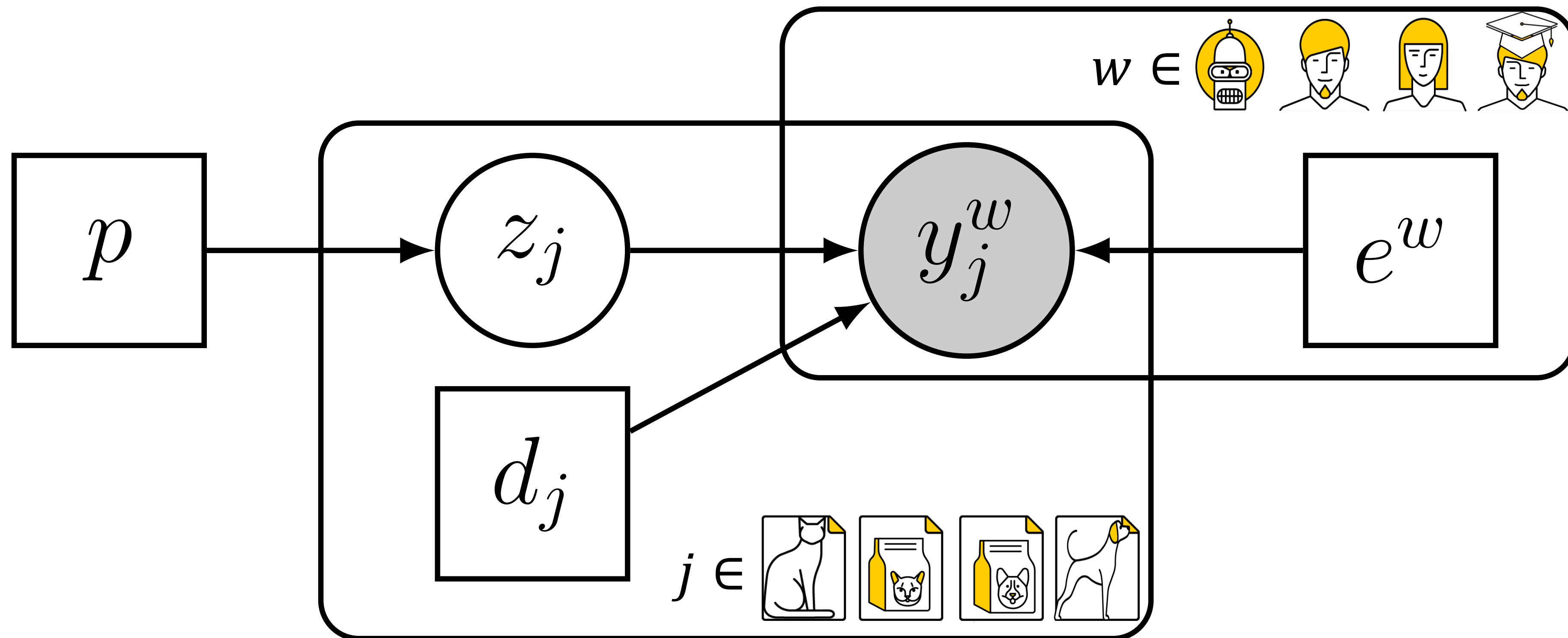
Advanced aggregation: performers and objects

› Parameterize expertise of performers by e^w

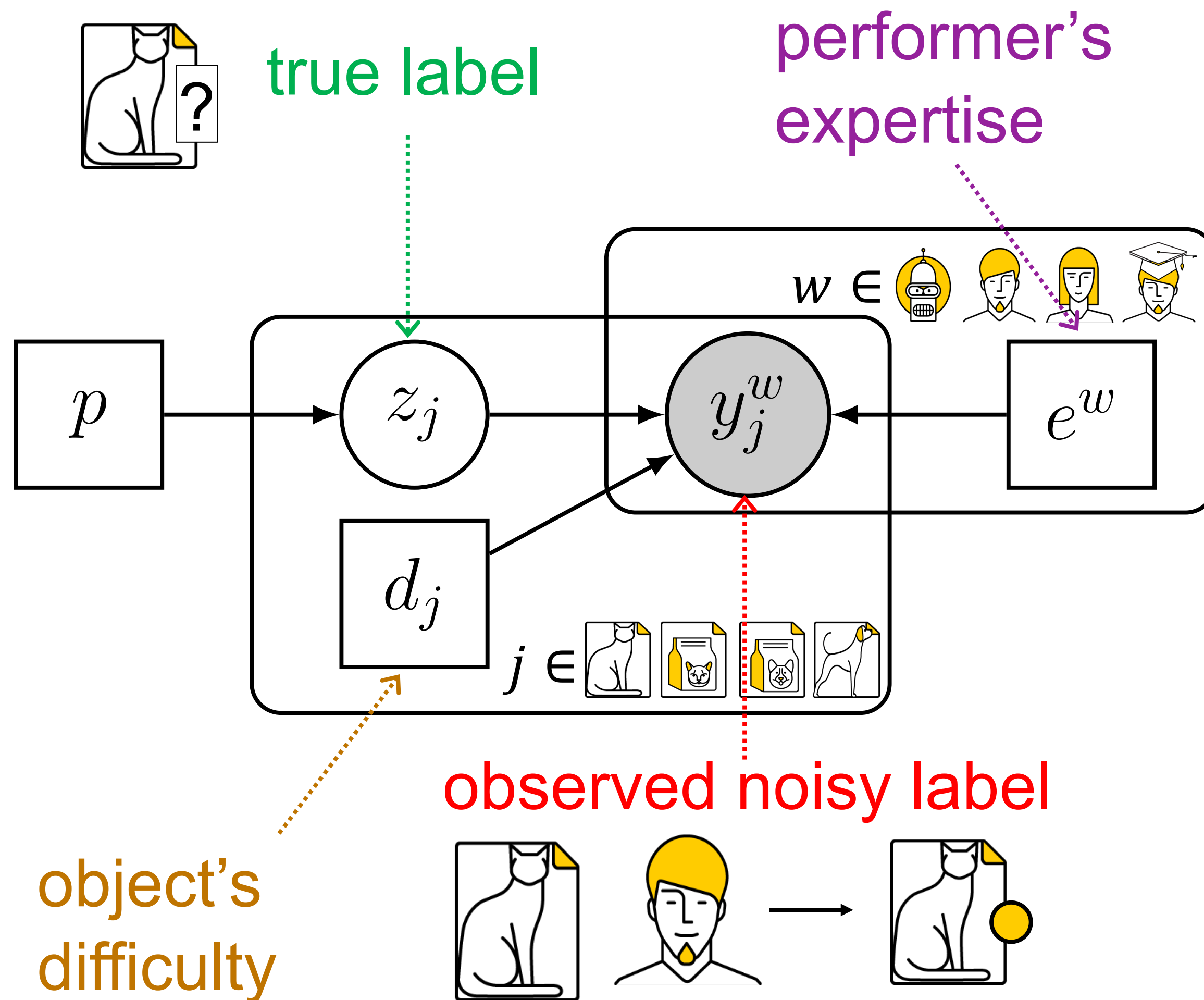
› Parameterize difficulty of objects by d_j



Advanced aggregation: latent label models

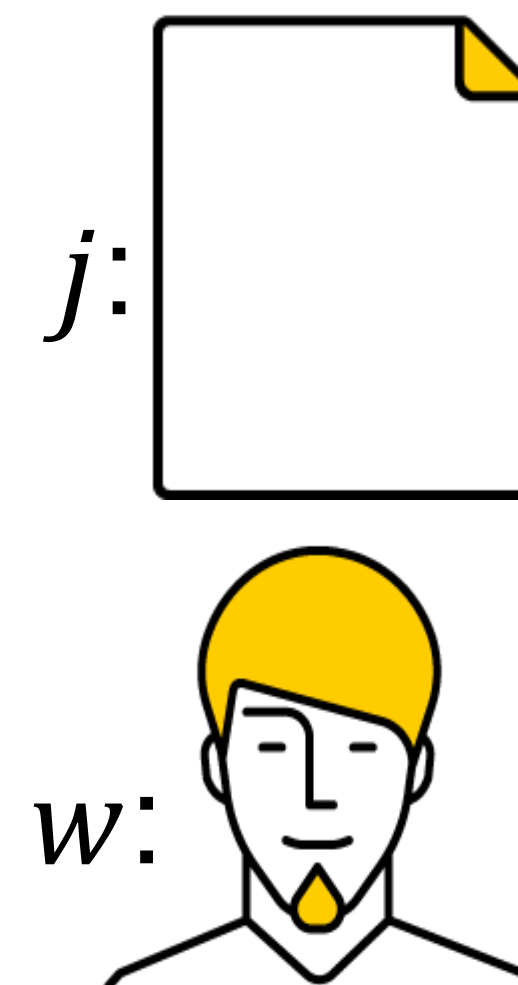


Latent label models: noisy label model



› A noisy label model $M_j^w = M(e^w, d_j)$ is a matrix of size $K \times K$ with elements

$$M_j^w[c, k] = \Pr(Y_j^w = k \mid Z_j = c)$$

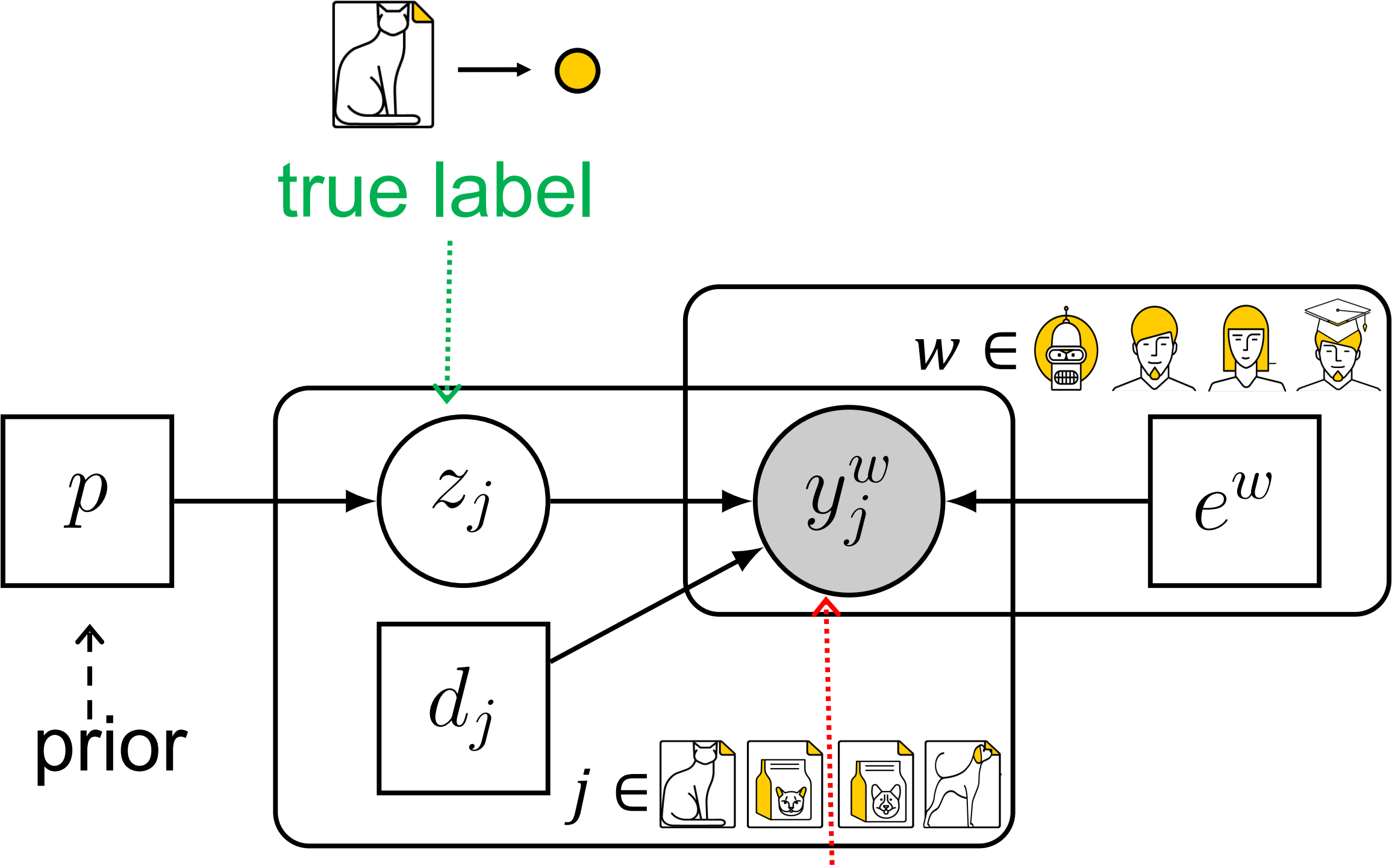


M_j^w :

Noisy \ True	●	▲	■
●	q_{11}	q_{12}	q_{13}
▲	q_{21}	q_{22}	q_{23}
■	q_{31}	q_{32}	q_{33}

$$q_{c1} + q_{c2} + q_{c3} = 1 \text{ for each } c$$

Latent label models: generative process



› Noisy labels generation:

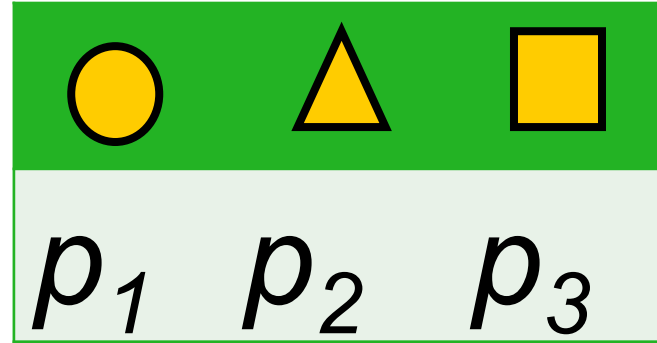
1. Sample z_j from a distribution

$$P_Z(p)$$

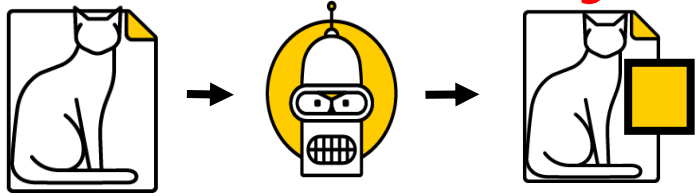
2. Sample y_j^w from a distribution

$$P_Y(M_j^w[z_j, \cdot])$$

In multiclassification, a standard choice for $P_Z(\cdot)$ and $P_Y(\cdot)$ is a Multinomial distribution $\text{Mult}(\cdot)$



observed noisy label



$z_j \backslash y_j^w$			
	q_{11}	q_{12}	q_{13}

Latent label models: parameters optimization

- › Assumption: y_j^w is cond. independent of everything else given z_j, d_j, e^w
- › The likelihood of \mathbf{y} and \mathbf{z} under the latent label model:

$$L\left(\underbrace{\{z_j\}_{j=1}^J}_{\text{latent true label}}, \underbrace{p, \{d_j\}_{j=1}^J, \{e^w\}_{w=1}^W}_{\text{latent parameters}}\right) = \prod_{j \in J} \underbrace{\sum_{z_j \in \{1, \dots, K\}} \Pr(z_j | p) \prod_{w \in W_j} \Pr(y_j^w | z_j, d_j, e^w)}_{\text{likelihood of noisy and true labels for object } j}$$

observed
noisy label

- › Estimate parameters and true labels by maximizing $L(\dots)$

Latent label models: EM algorithm

- › Maximization of the expectation of log-likelihood (LL), a lower bound on LL of \mathbf{y} and \mathbf{z}

$$\mathbb{E}_{\mathbf{z}} \log \Pr(\mathbf{y}, \mathbf{z}) = \sum_{j \in J} \sum_{z_j \in \{1, \dots, K\}} \Pr(z_j | p) \log \prod_{w \in W_j} \Pr(z_j | p) \Pr(y_j^w | z_j, d_j, e^w)$$

- › **E-step:** Use Bayes' theorem for posterior distribution of $\hat{\mathbf{z}}$ given $p, \mathbf{d}, \mathbf{e}$:

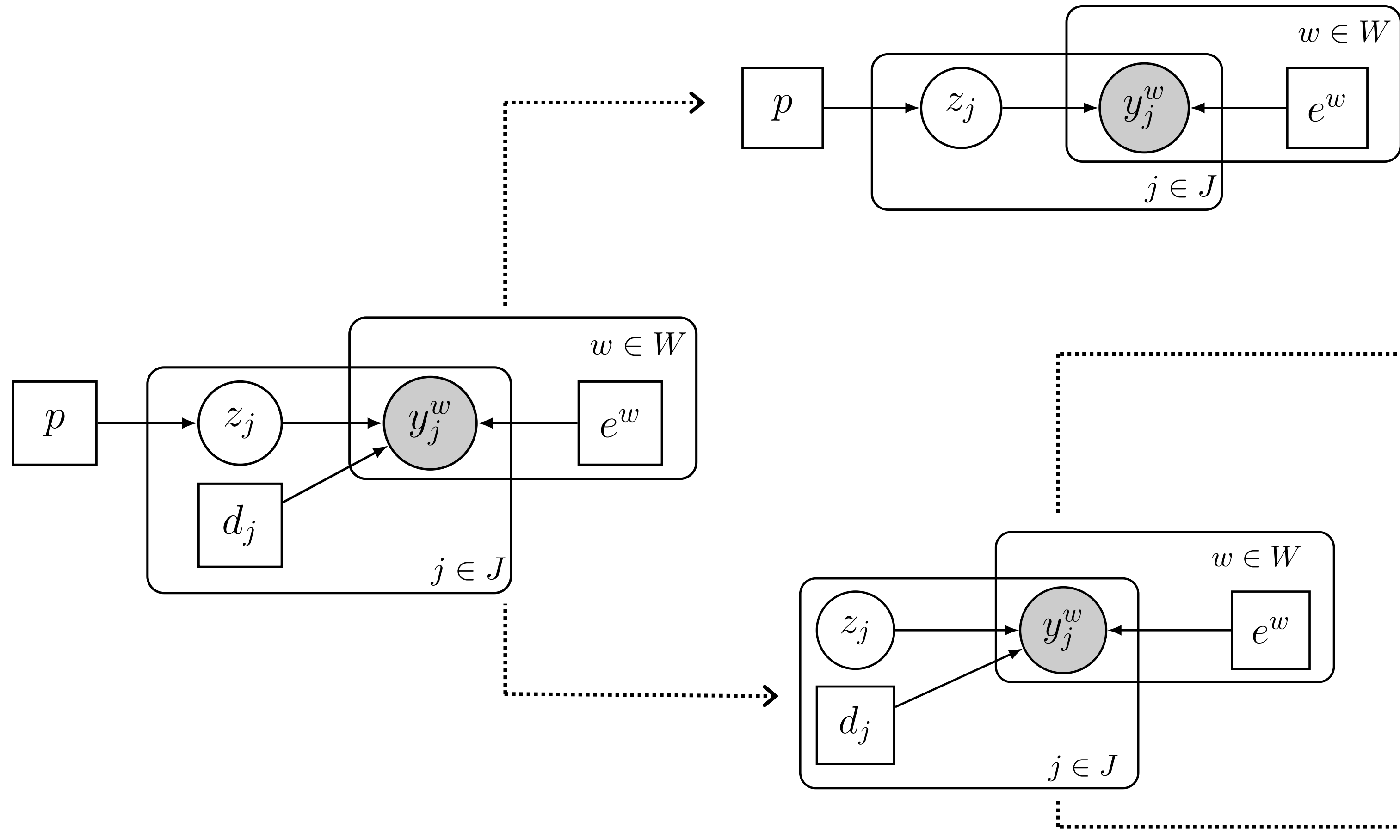
$$\hat{z}_j[c] = \Pr(Z_j = c | \mathbf{y}, p, \mathbf{d}, \mathbf{e}) \propto \Pr(Z_j = c | p) \prod_{w \in W_j} \Pr(y_j^w | Z_j = c, d_j, e^w)$$

- › **M-step:** Maximize the expectation of LL with respect to the posterior distribution of $\hat{\mathbf{z}}$:

$$(p, \mathbf{d}, \mathbf{e}) = \operatorname{argmax} \mathbb{E}_{\hat{\mathbf{z}}} \log \Pr(z_j | p) \prod_{w \in W_j} \Pr(y_j^w | z_j, d_j, e^w)$$

- Analytical solutions
- Gradient descent

Latent label model (LLM): special cases



Dawid and Skene model (DS):

- › categories are different
- › objects are similar
- › performers are different

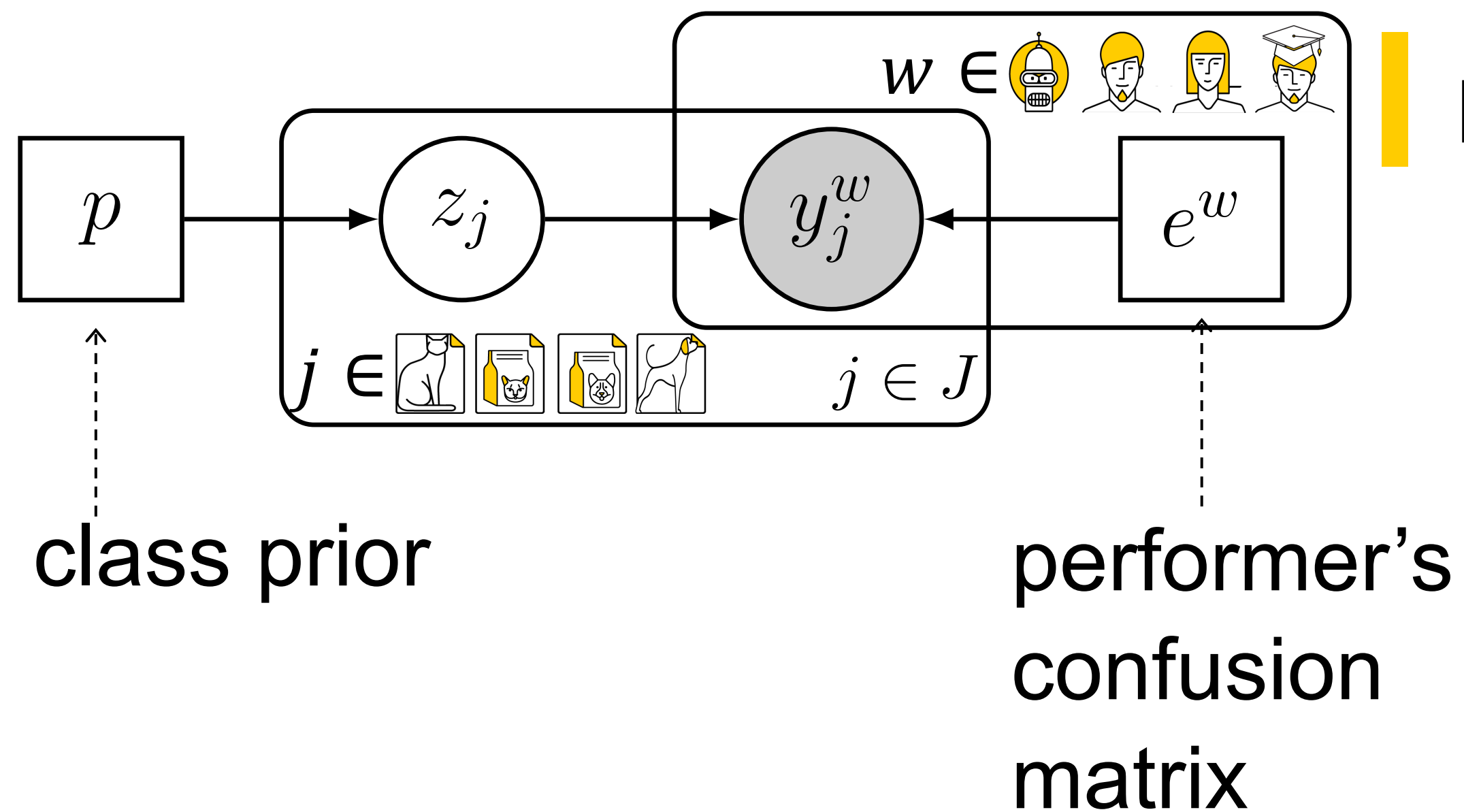
Generative model of labels, abilities, and difficulties (GLAD):

- › categories are similar
- › objects are different
- › performers are different

Minimax conditional entropy model (MMCE):

- › categories are different
- › objects are different
- › performers are different

Dawid and Skene model (DS)



LLM with parameters:

- > p – vector of length K : $p[i] = \Pr(\mathbf{Z} = c)$
- > e^w – matrix of size $K \times K$:

$$e^w[c, k] = \Pr(\mathbf{Y}^w = k | \mathbf{Z} = c)$$

$\mathbf{z} \backslash \mathbf{y}_w$	○	△	□
○	■	—	■
△	—	■	■
□	■	■	■

DS: parameters optimization

› **E-step:**

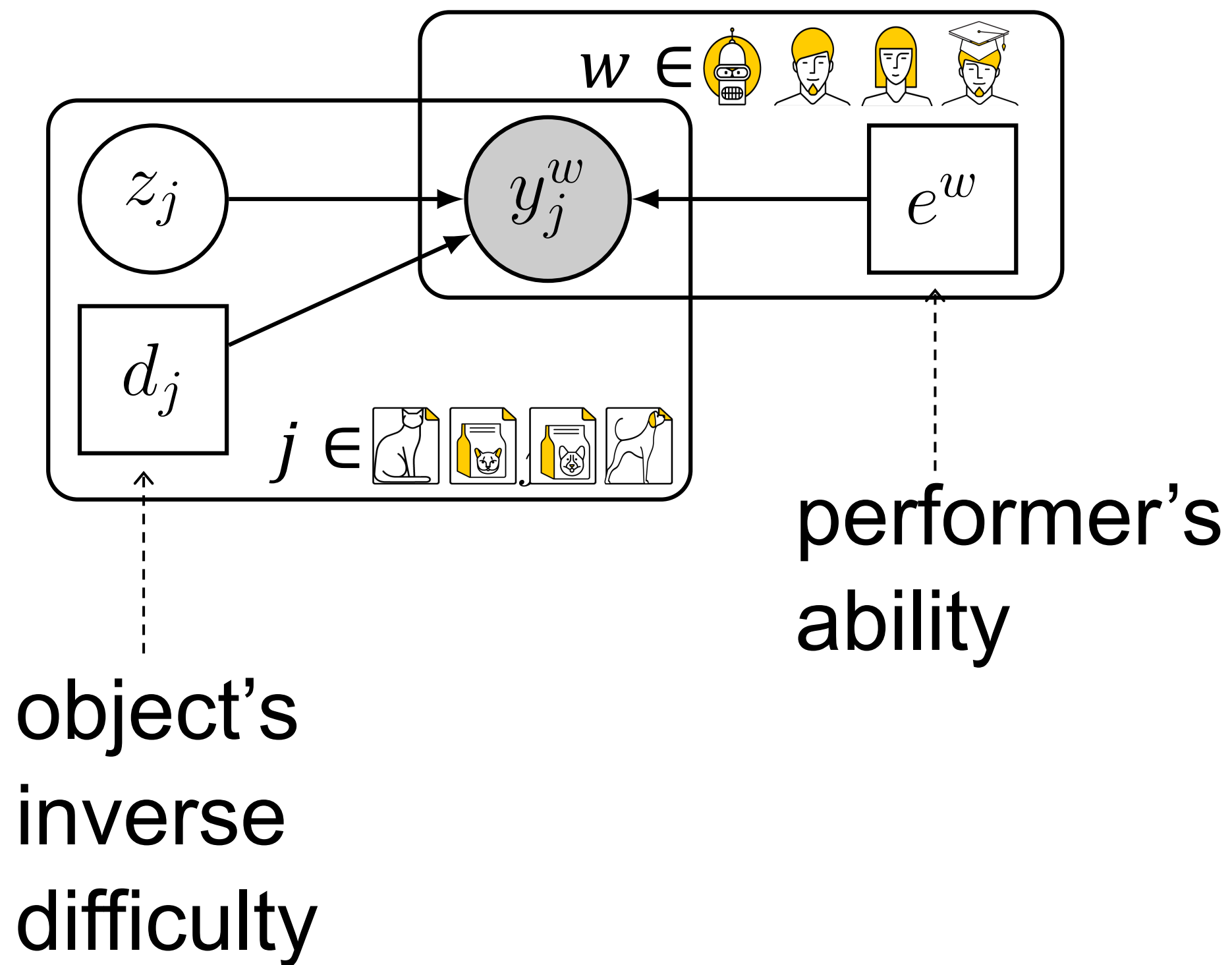
$$\hat{z}_j[c] = \frac{p[c] \prod_{w \in W_j} e^w[c, y_j^w]}{\sum_k p[k] \prod_{w \in W_j} e^w[k, y_j^w]}, \quad c = 1, \dots, K$$

› **M-step:** Analytical solution

$$e^w[c, k] = \frac{\sum_{j \in J} \hat{z}_j[c] \delta(y_j^w = k)}{\sum_{q=1}^K \sum_{j \in J} \hat{z}_j[c] \delta(y_j^w = q)}, \quad k, c = 1, \dots, K$$

$$p[c] = \frac{\sum_{j \in J} \hat{z}_j[c]}{J}, \quad c = 1, \dots, K$$

Generative model of Labels, Abilities, and Difficulties (GLAD)



LLM with parameters:

- › scalar $d_j \in (0, \infty)$
- › scalar $e^w \in (-\infty, \infty)$
- › Model:

$$\Pr(Y_j^w = k | Z_j = c) = \begin{cases} a(w, j), & c = k \\ \frac{1 - a(w, j)}{K - 1}, & c \neq k \end{cases}$$

$$\text{where } a(w, j) = \frac{1}{1 + \exp(-e^w d_j)}$$

GLAD: parameters optimization

- › Let $a(w, j) = \frac{1}{1 + \exp(-e^w d_j)}$ and $P(z_j)$ be a predefined prior (e.g., $P(z_j) = 1/K$)
- › **E-step:**

$$\hat{z}_j [c] \propto P(Z_j = c) \prod_{w \in W_j} a(w, j)^{\delta(y_j^w = c)} \left(\frac{1 - a(w, j)}{K - 1} \right)^{\delta(y_j^w \neq c)}, \quad c = 1, \dots, K$$

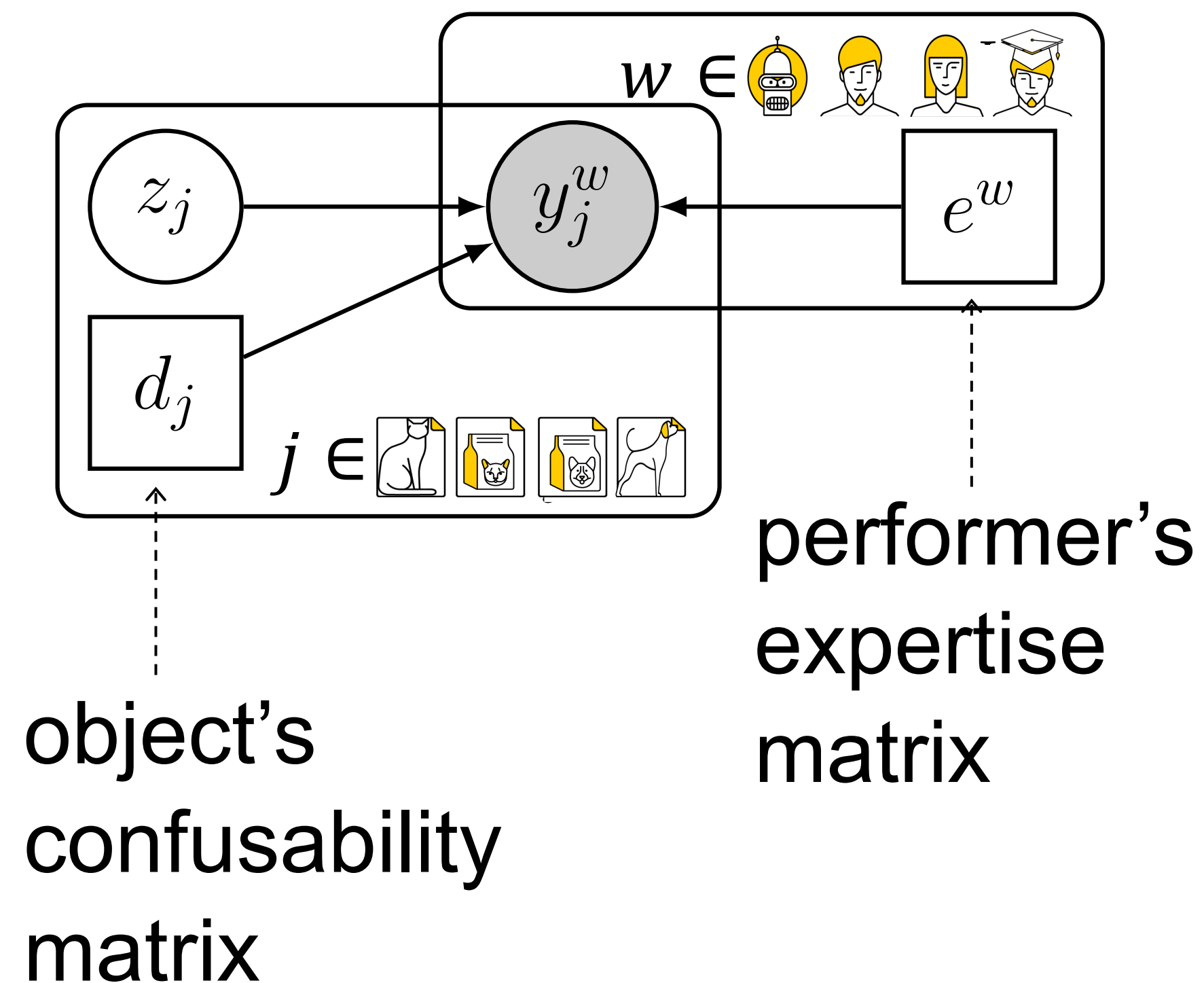
- › **M-step:** estimate (\mathbf{d}, \mathbf{e}) for given $\hat{\mathbf{z}}$ using gradient descent

$$(\mathbf{d}^t, \mathbf{e}^t) = \operatorname{argmax} \sum_{j \in J} \left[\mathbb{E}_{\hat{z}_j} \log P(z_j) + \sum_{w \in W_j} \mathbb{E}_{\hat{z}_j} \log \Pr(y_j^w | z_j) \right]$$

MiniMax Conditional Entropy model (MMCE)

- Find parameters that minimize the maximum conditional entropy of observed labels:

$$\min_Q \max_P - \sum_{j \in J} \sum_{c \in \{1, \dots, K\}} Q(Z_j = c) \sum_{w \in W} \sum_{k \in \{1, \dots, K\}} P(Y_j^w = k | Z_j = c) \log P(Y_j^w = k | Z_j = c)$$

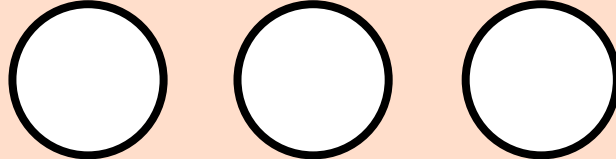
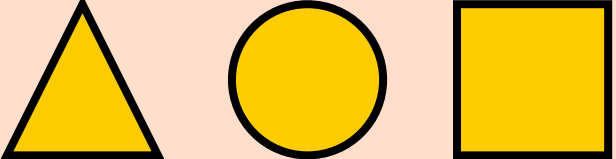
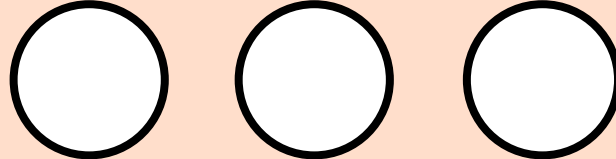
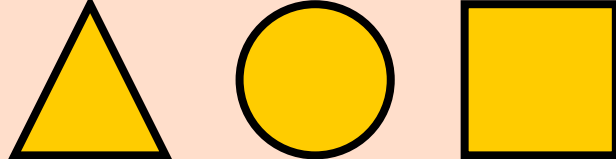
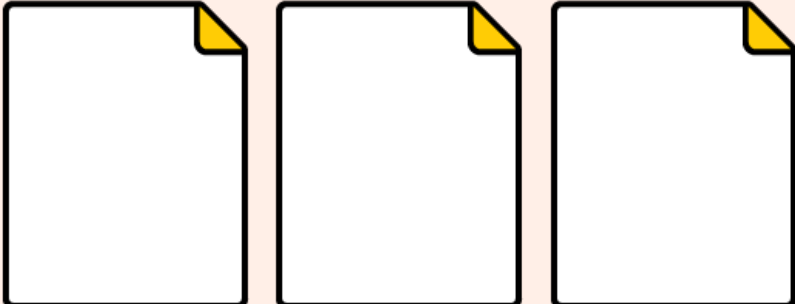
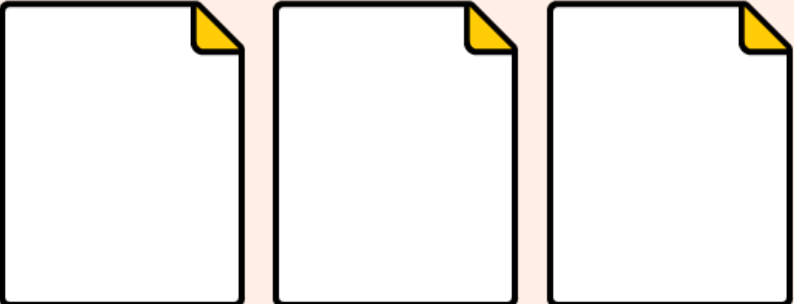

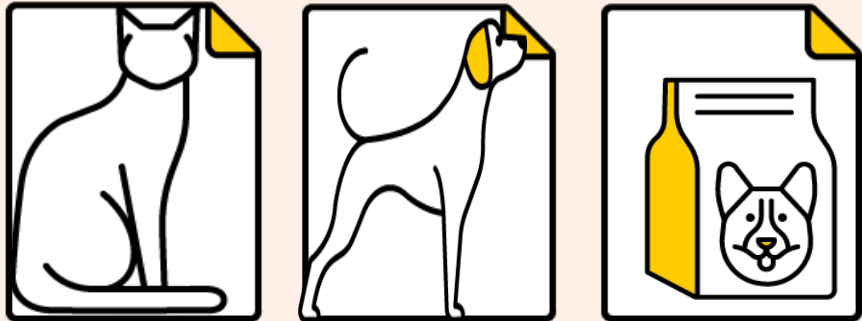






LLM with parameters:

- d_j — matrix of size $K \times K$
- e^w — matrix of size $K \times K$
- Noisy label model:

$$\Pr(Y_j^w = k | Z_j = c) = \exp(d_j[c, k] + e^w[c, k])$$

Summary of aggregation methods

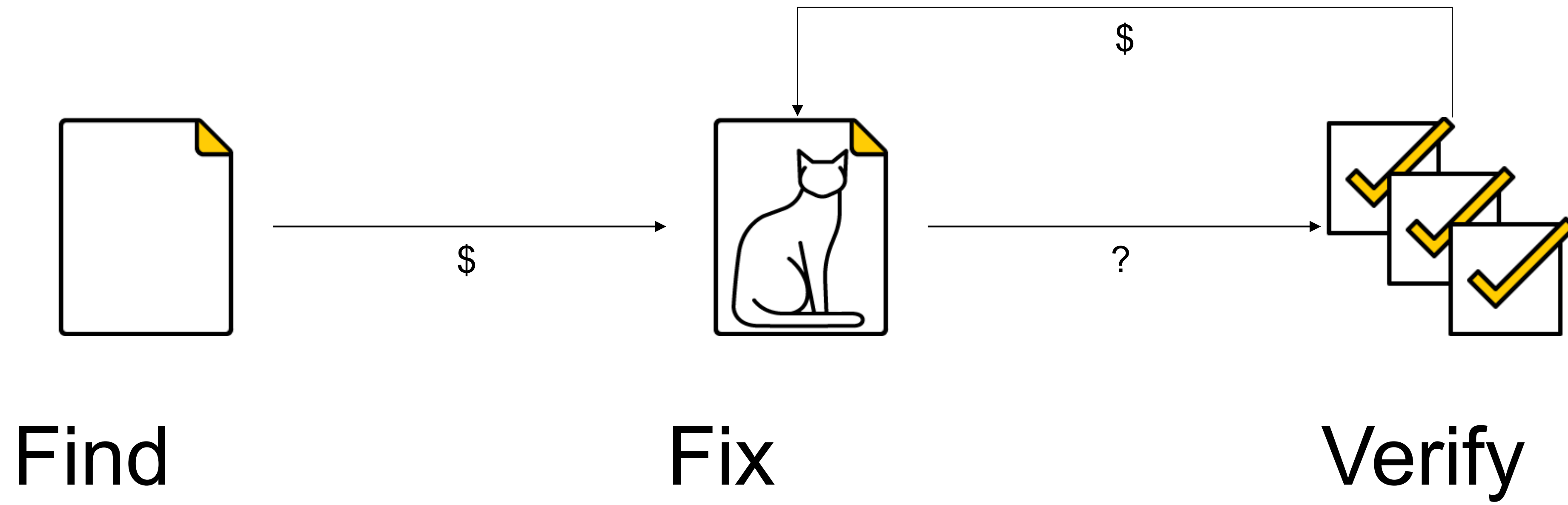
	MV	DS	GLAD	MMCE
Categories (K)				
Objects (J)				
Performers (W)				
Number of parameters	0	$WK^2 + K$	$W + J$	$(W + J)K^2$

Text Aggregation

Text Aggregation

- So far we discussed how to aggregate *categorical* responses
- In NLP we often work with **textual data**, i.e., with *sequences*
- How can we solve tasks with "unknown" responses?

Crowdsourced Copy-Editing: Soylent



Automatic Text Aggregation

- Post-acceptance is a universal technique for open-ended tasks
- However, it adds additional (slight) complexity to the pipeline
- Can we aggregate texts without human intervention?
- We would like to minimize **Word Error Rate** (WER) computed as a function of the number of **C**orrect, **S**ubstitution, **D**eletion, **I**nsertion items:

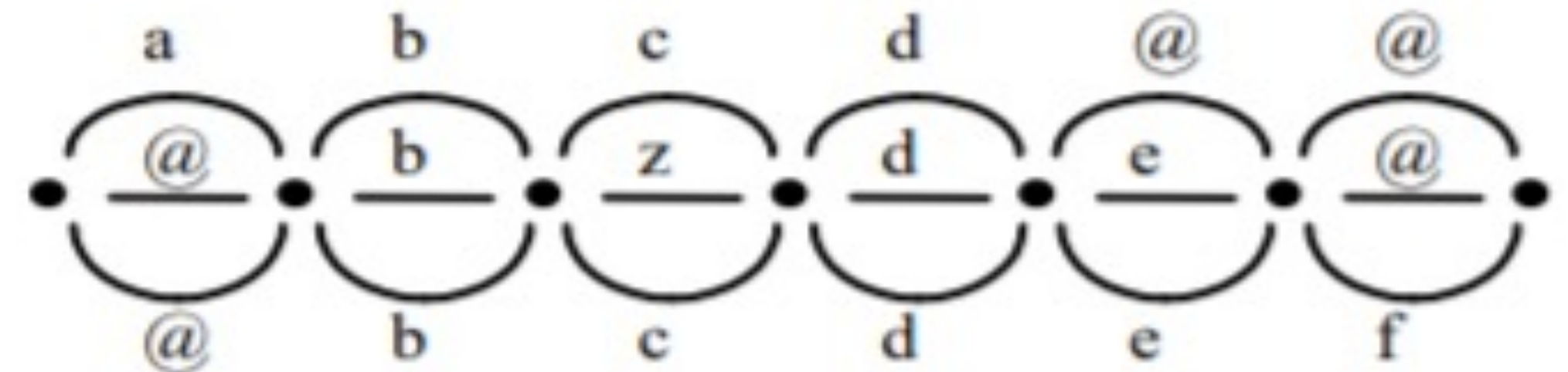
$$\text{WER} = \frac{S + D + I}{C + S + D}$$

Automatic Text Aggregation: ROVER

An efficient method for long sequences:

- **Input:** a b c d; b z d e; b c d e f

- **Word Transition Network:**
(words with highest scores are chosen)



- **Result:** b c d e

Automatic Text Aggregation: HRRASA

- Obtain the sequence embeddings with BERT, RoBERTa, etc.
- Choose the response that is the closes to the embedding $e()$ of the estimated response (a_j^w) provided by performer w for task j

$$\beta_w = \frac{X_{(\frac{\alpha}{2}, |V_w|)}^2}{\Sigma(e(a_j^w) - \hat{e}_j)^2} \quad \hat{e}_j = \frac{\Sigma \beta_w e(a_j^w)}{\Sigma \beta_w}$$

$$s_j^w = \beta_w \cdot \exp\left(-\frac{\|e_j^w - \hat{e}_j\|^2}{\|e_j^w\|^2 \|\hat{e}_j\|^2}\right) + \gamma_j^w$$

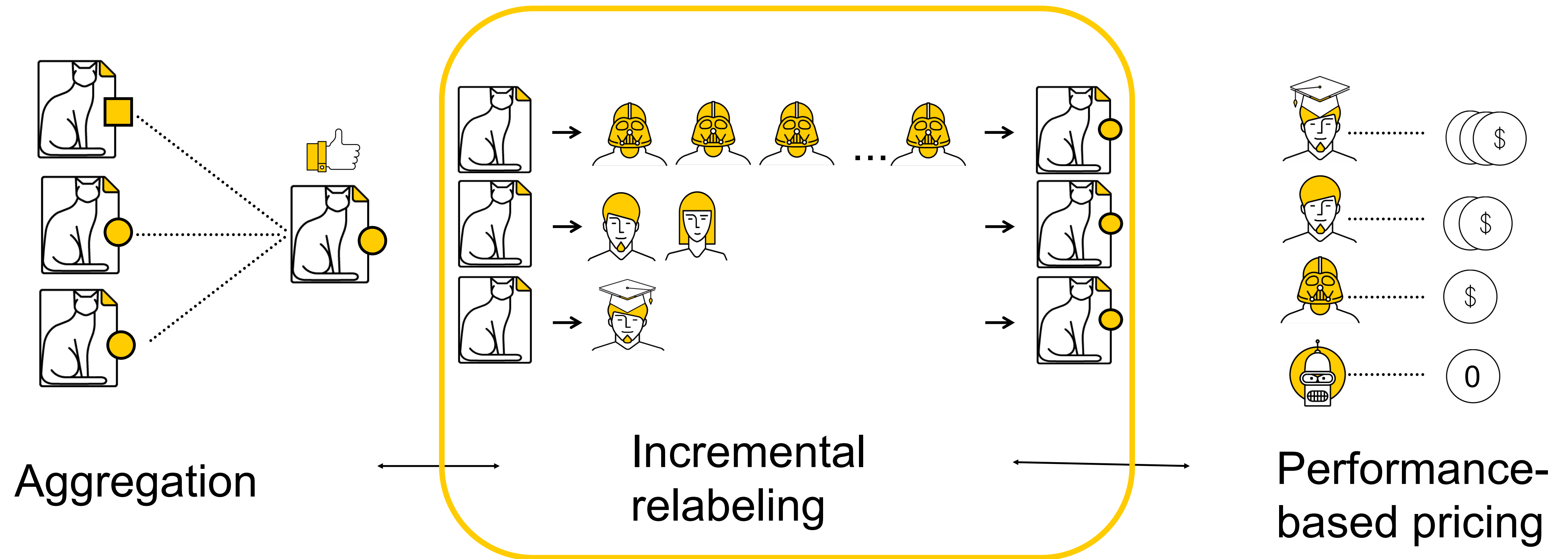
- Parameters are estimated step-by-step

Crowd-Kit, a General-Purpose Toolkit

<https://pypi.org/project/crowd-kit/>

Crowd-Kit allows aggregating answers of many kinds, including categorical, sequential, and graphical, using the same API.

Key components of labeling with crowds



Incremental relabeling

aka dynamic overlap

Pool settings: dynamic overlap

Quality control

Add rules to get more accurate responses.
All rules work independently.

NON-AUTOMATIC ACCEPTANCE ?

No

REVIEW PERIOD IN DAYS

CAPTCHA FREQUENCY ?

None

Add Quality Control Rule

Overlap

Specify how many performers you want to complete each task in the pool.

OVERLAP ?

DYNAMIC OVERLAP ?

Off

Speed/quality ratio

Specify additional conditions for selecting performers by their rating in Toloka.
This will improve quality, but may reduce the speed of task completion because there will be fewer performers available for completing tasks. [Learn more](#)

Top %

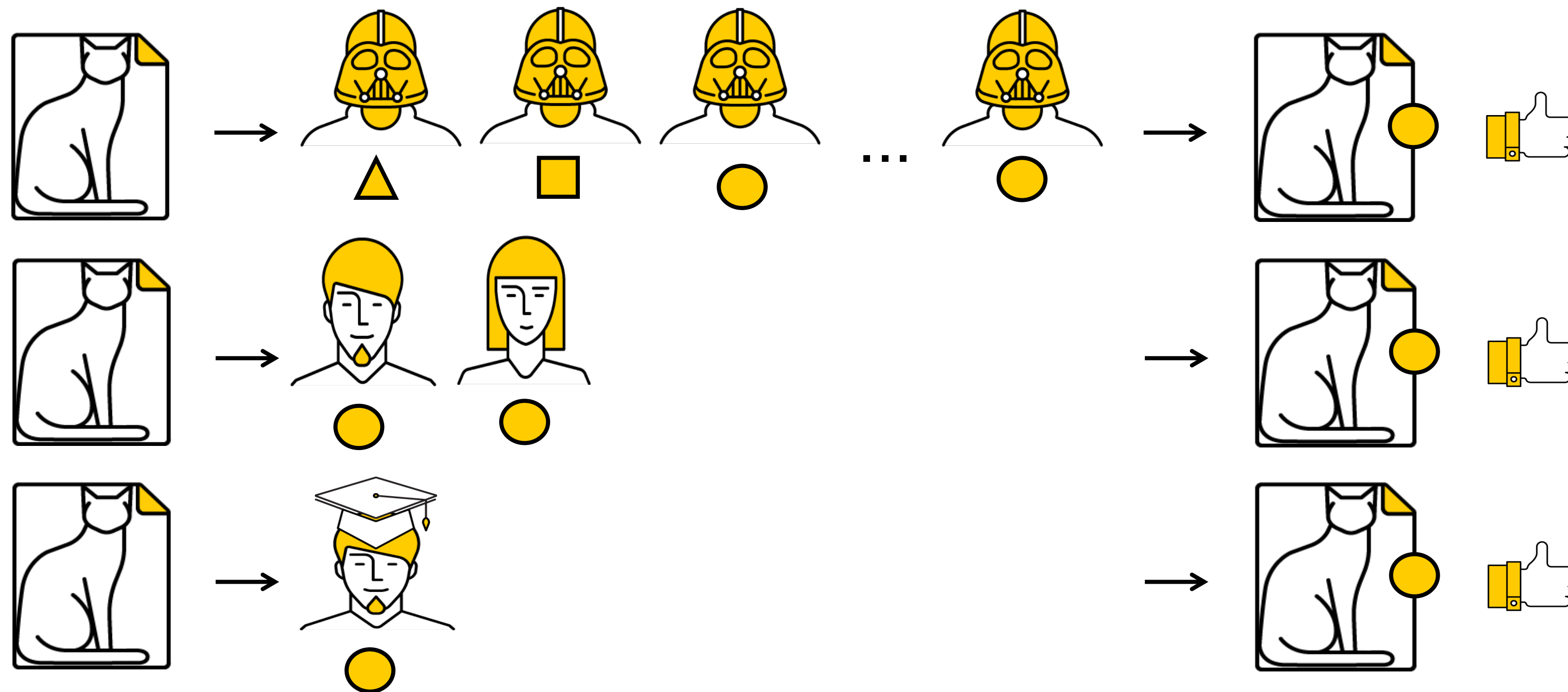
Online

Time

Specify the percentage of top-rated active users who can access tasks in the pool.

Incremental relabeling problem

Obtain aggregated labels of a desired level of quality using a fewer number of noisy labels



Incremental relabeling scheme (IRL)

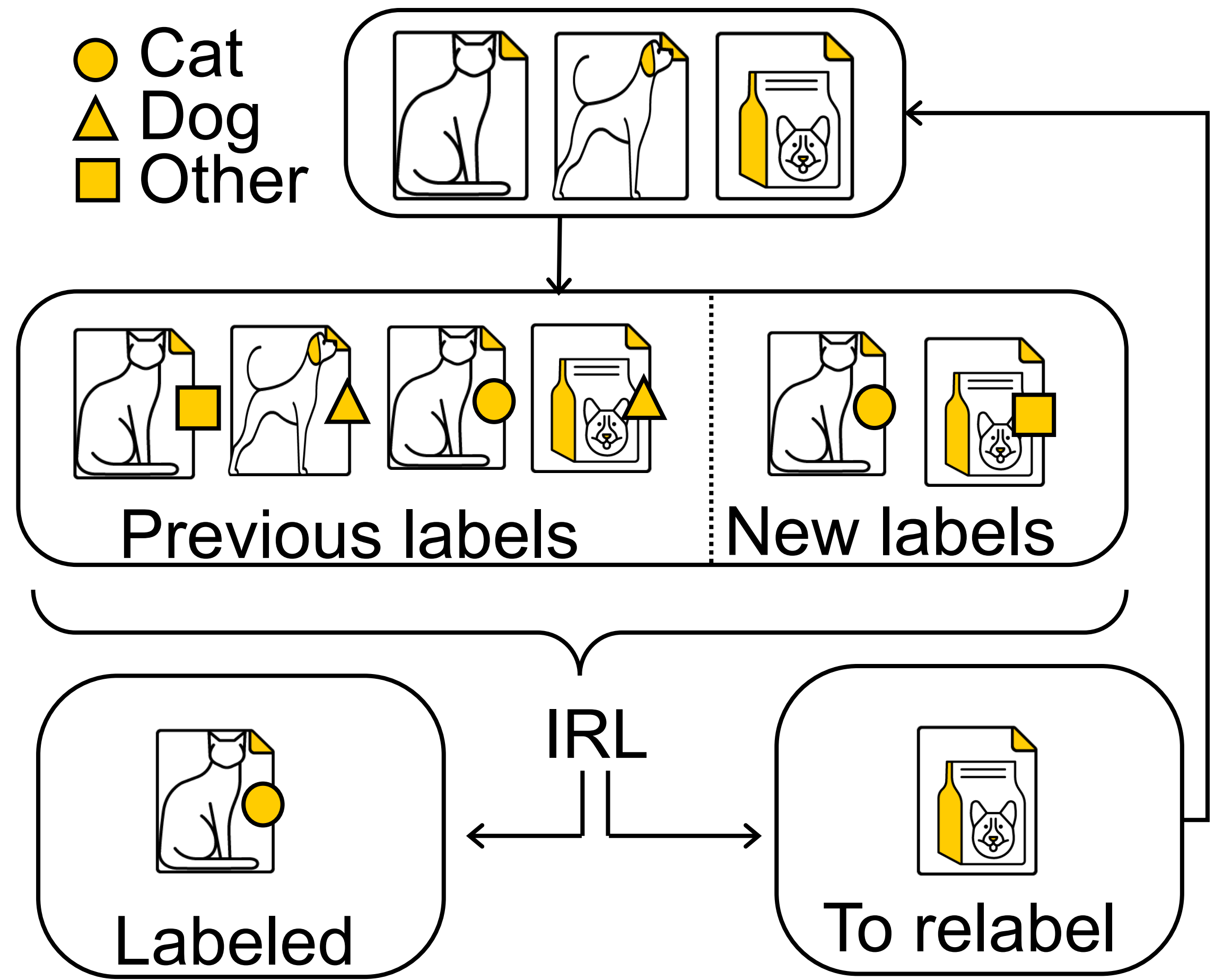
Request a label for each object

In real time IRL algorithm receives:
(1) previously accumulated labels
(2) new labels

Decides:

(1) which objects are labeled
(2) which objects to relabel

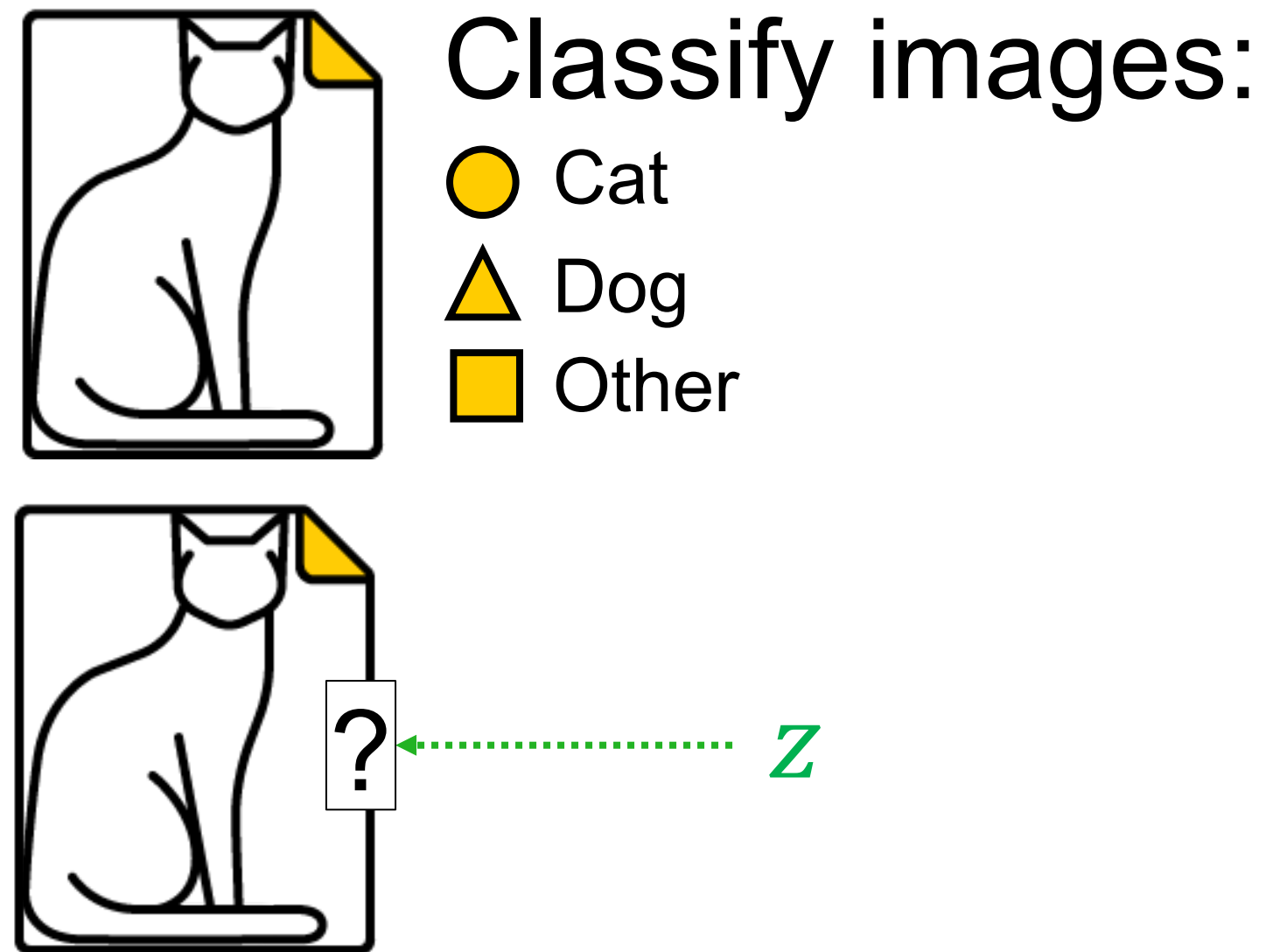
Repeat until all tasks are labeled



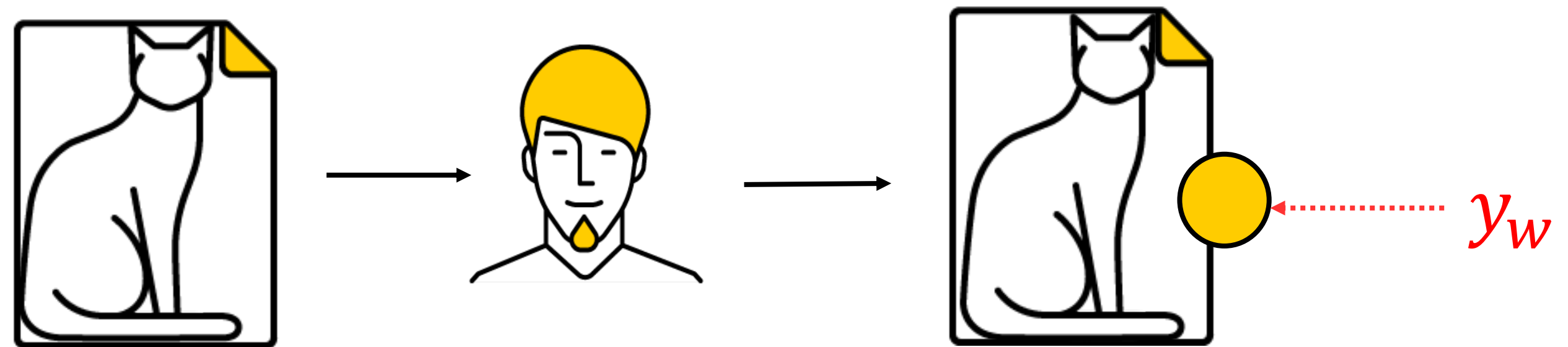
Notations

› Consider one object

› $z \in \{1, \dots, K\}$ - latent true label



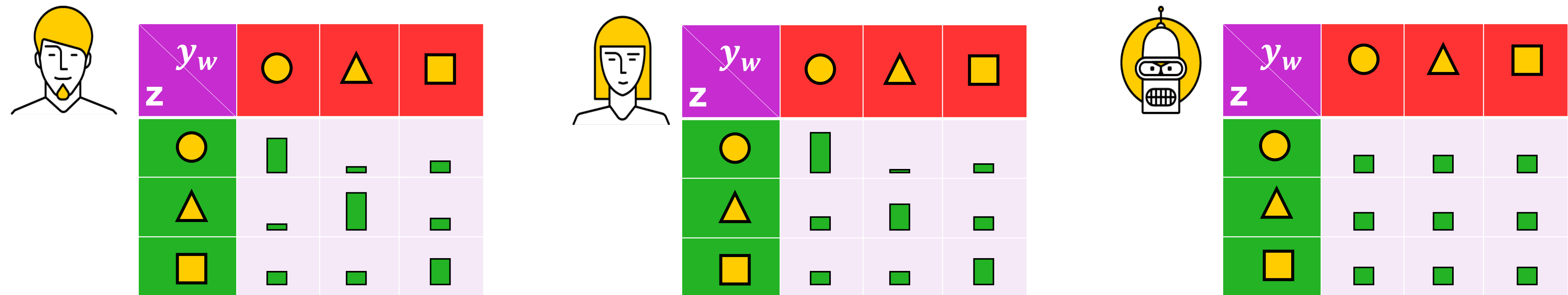
› $y_w \in \{1, \dots, K\}$ - observed noisy label from performer w :



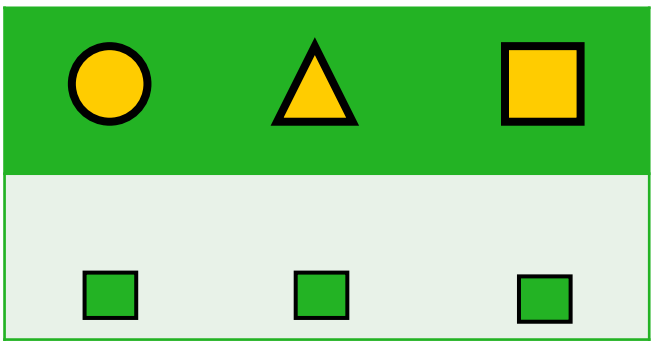
Notations

› Noisy label model for performer w :

$$M_w \in [0,1]^{K \times K}: \Pr(Y_w = k | Z = c) = M_w[c, k]$$



› Prior distribution: $\Pr(Z = k) = p_k$

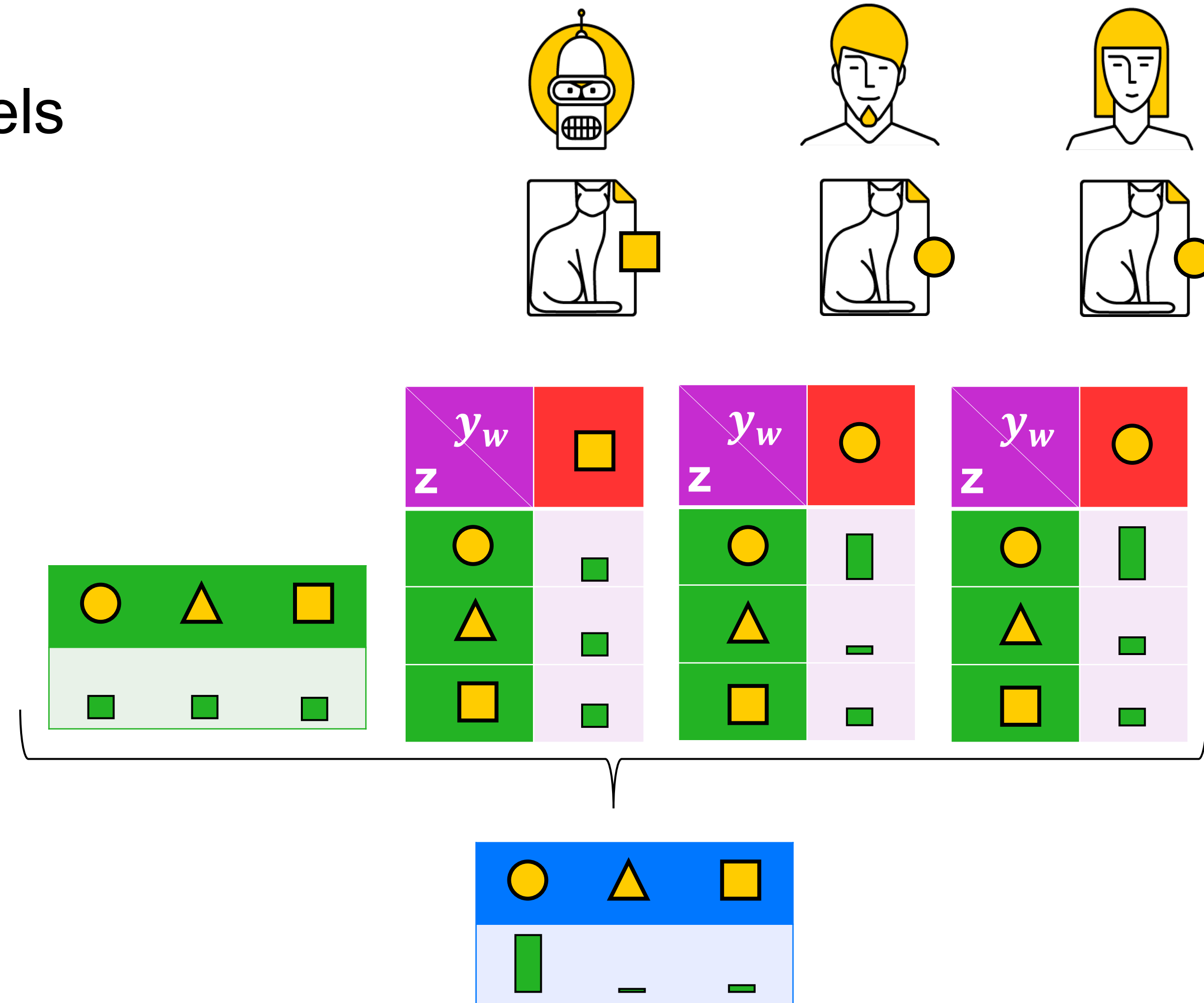


Posterior distribution

- › $\{y_{w_1}, \dots, y_{w_n}\}$ - accumulated noisy labels for the object

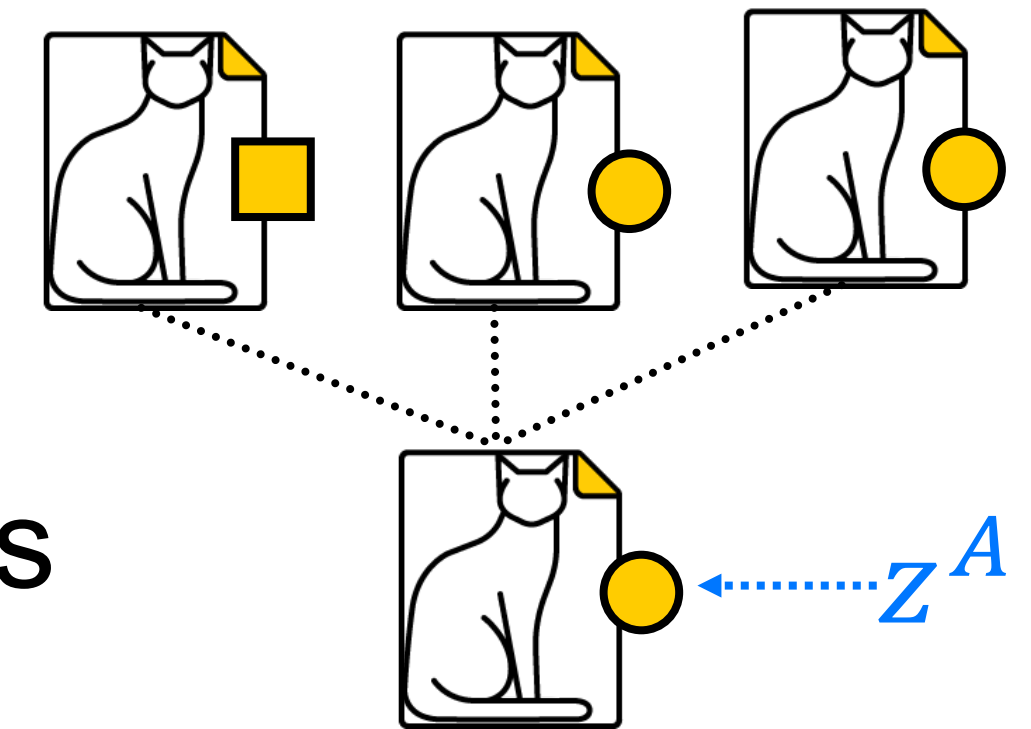
- › Using Bayes rule:

$$\begin{aligned} & \Pr(Z = k | \{y_{w_1}, \dots, y_{w_n}\}) \\ &= \frac{\Pr(Z = k) \Pr(\{y_{w_1}, \dots, y_{w_n}\} | Z = k)}{\Pr(\{y_{w_1}, \dots, y_{w_n}\})} \\ &= \frac{p_k \prod_{i=1}^n M_{w_i}[k, y_{w_i}]}{\sum_{t=1}^K p_t \prod_{i=1}^n M_{w_i}[t, y_{w_i}]} \end{aligned}$$



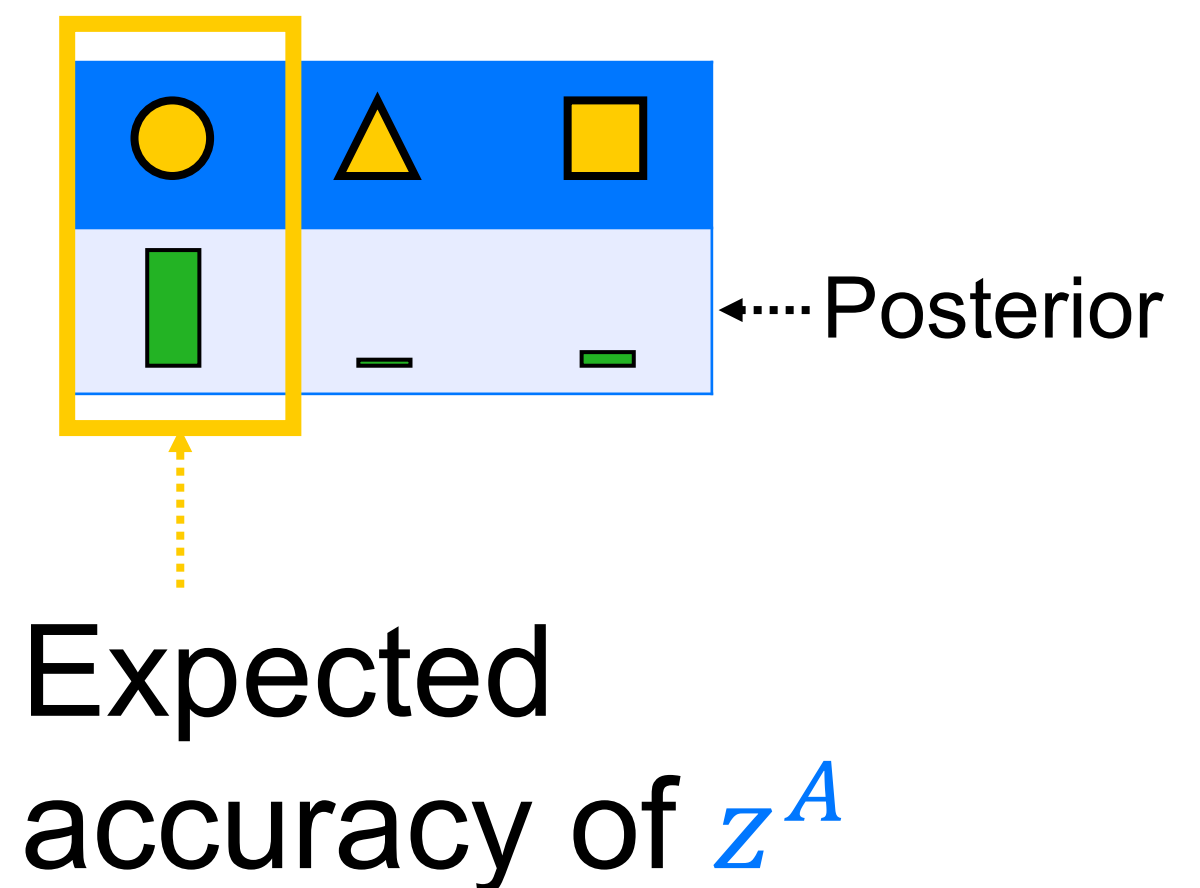
Expected accuracy of aggregated labels

- › Let A be an aggregation model, e.g. MV, DS, GLAD,...
- › Denote aggregated label $z^A = A(\{y_{w_1}, \dots, y_{w_n}\})$
- › Expected accuracy of aggregated labels given noisy labels is

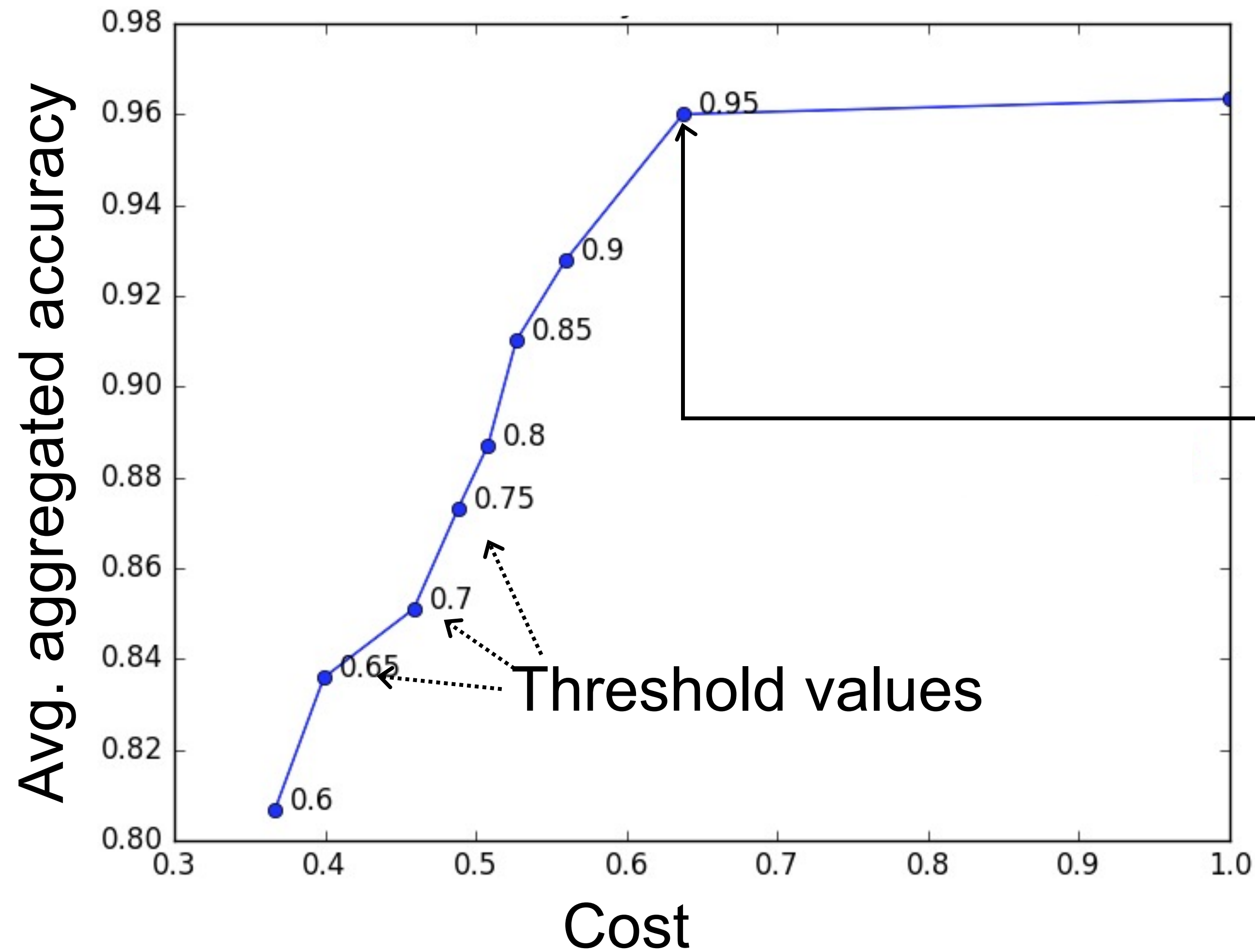


$$E(\delta(z = z^A) | \{y_{w_1}, \dots, y_{w_n}\}) = \Pr(z = z^A | \{y_{w_1}, \dots, y_{w_n}\})$$

- › Stop labeling if $E(\delta(z = z^A) | \{y_{w_1}, \dots, y_{w_n}\}) \geq c$
- \vdots
 parameter



Threshold in IRL: cost – accuracy trade-off

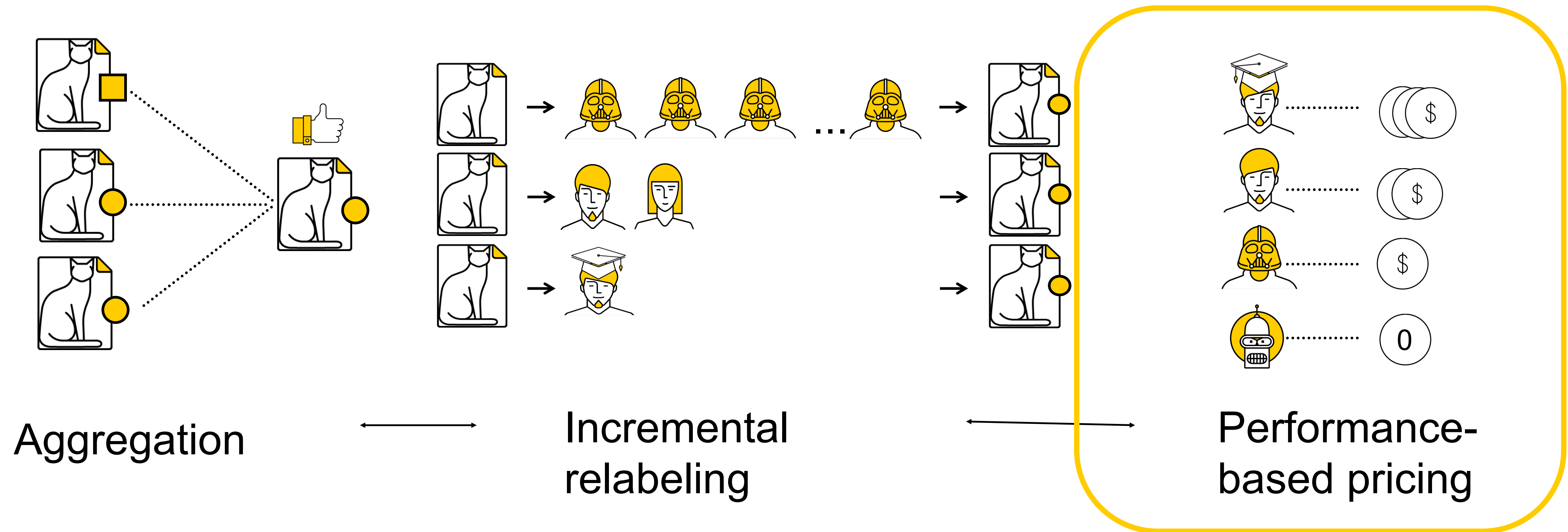


Optimal threshold $c = 0.95$

A higher c does not increase accuracy

Saving $\approx 35\%$ of noisy labels

Key components of labeling with crowds



Performance-based pricing

aka dynamic pricing

Pool settings: dynamic pricing

POOL NAME (VISIBLE ONLY TO YOU) ?

Are there traffic lights in the picture? ✕

☒ Use project description

PUBLIC DESCRIPTION ?

Add a private description

Price per task suite

You can add one or more tasks to the page. Enter the total price for all tasks on the page.

PRICE IN US DOLLARS ?

0.07

FEE ?

+ Dynamic pricing

Performers

Copy settings from...

Filter performers who can access the task.
Toloka has users from different countries, so don't forget to filter by language and region. [Learn more](#)

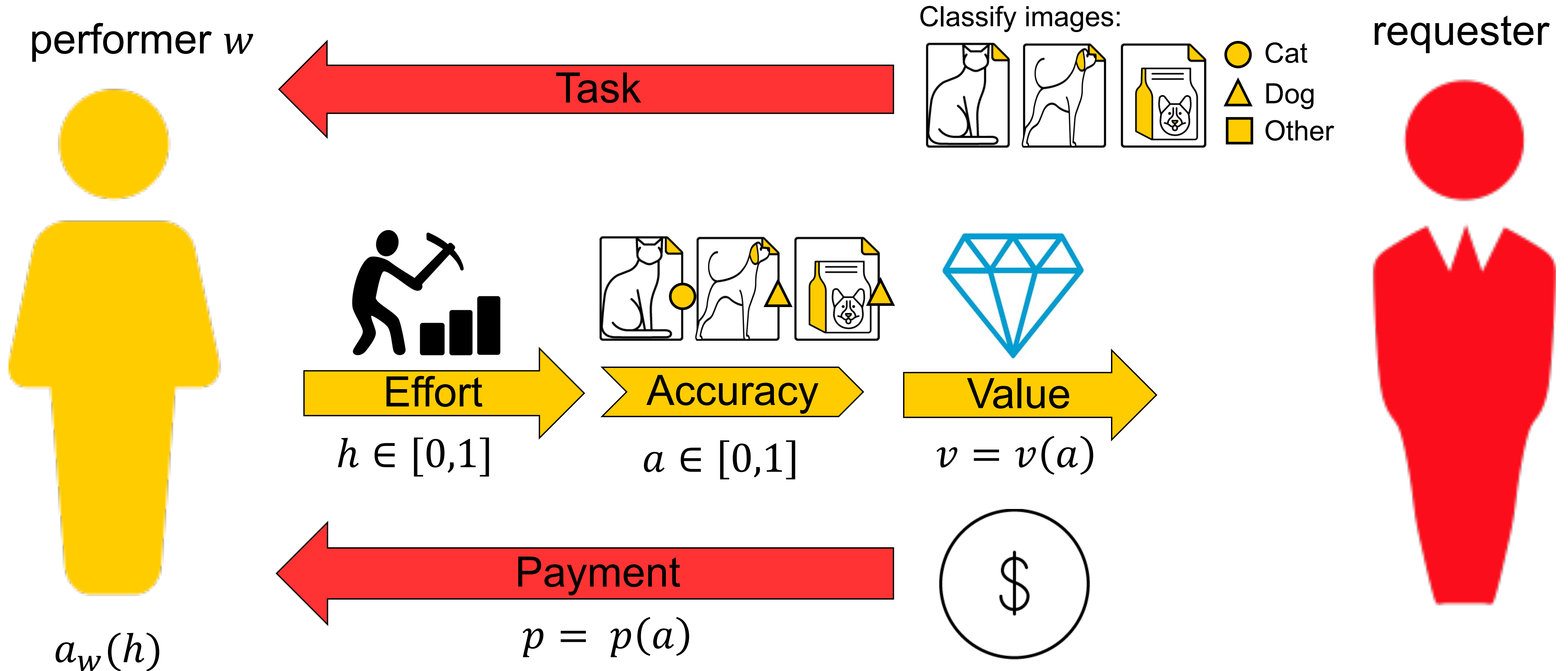
ADULT CONTENT ?

Yes ☐

Add filter ▼

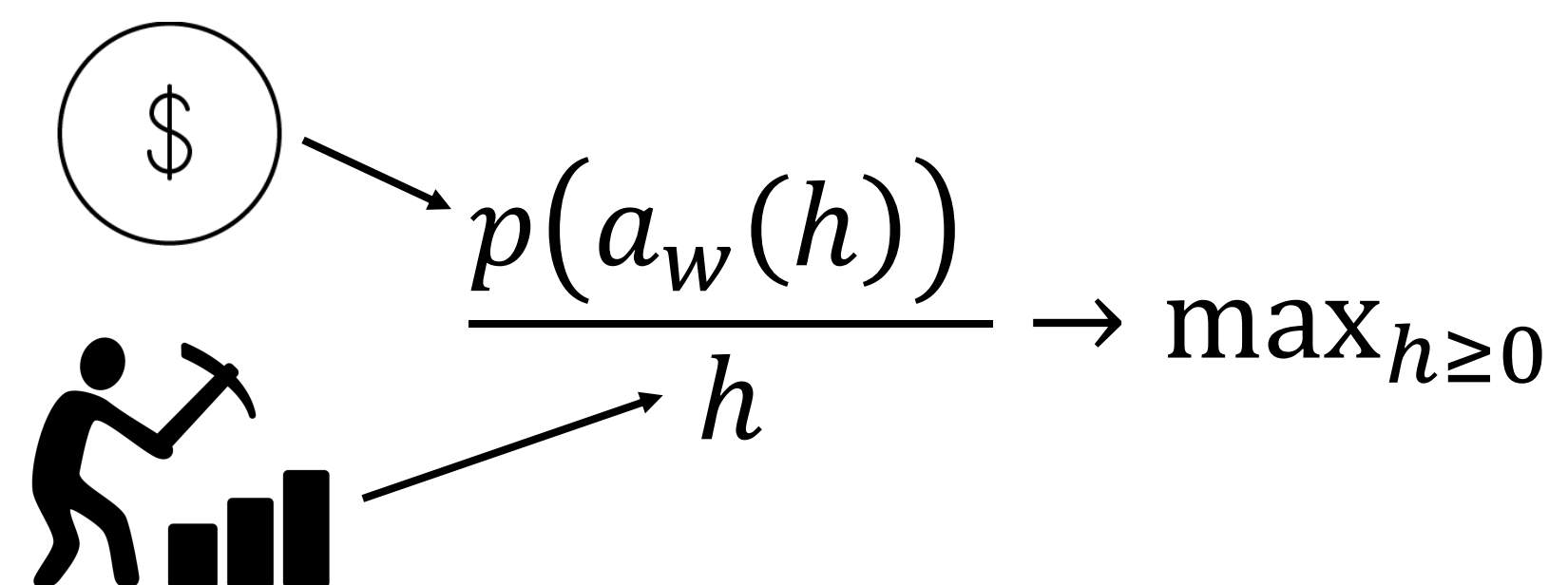
Create skill

Labeling as a game: notation

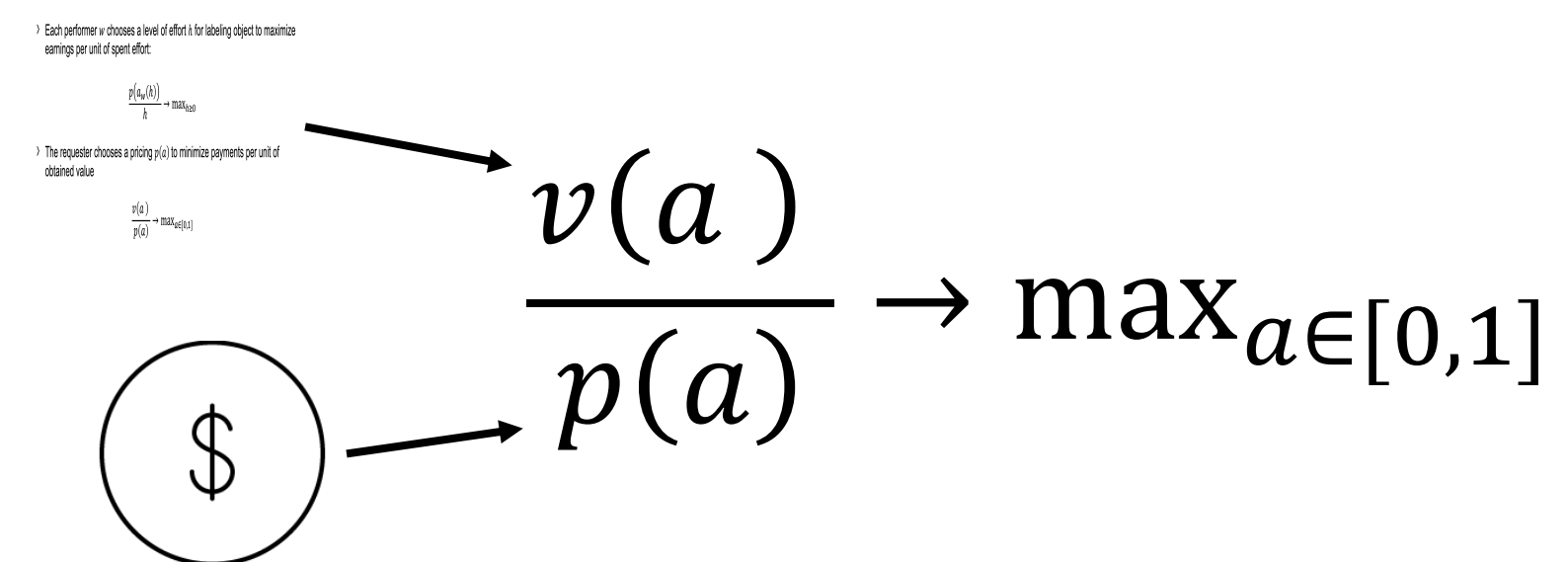


Labeling as a game: formalization

- Each performer w chooses a level of effort h for labeling object to maximize earnings per unit of spent effort:




- The requester chooses a pricing $p(a)$ to minimize payments per unit of obtained value



Labeling as a game: incentive compatible pricing

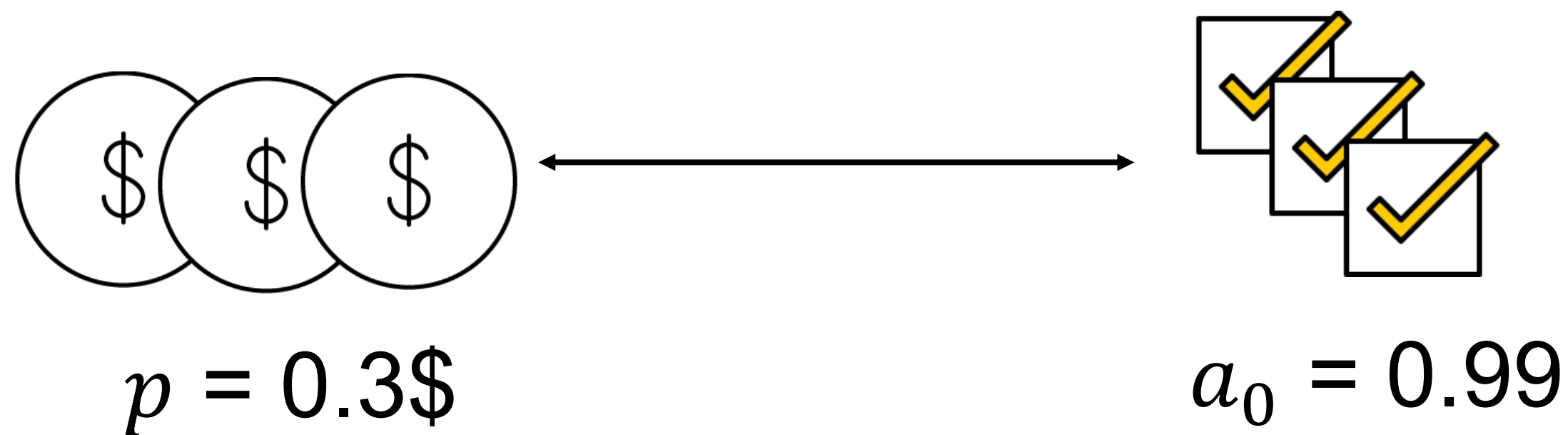
- › Assume $a_w(h)$ is a linear function of h :

$$\begin{array}{c} \nearrow a_w(h) = c_1 h + c_0 \\ \text{Accuracy} \end{array} \quad \begin{array}{c} \nwarrow \\ \text{Miner} \end{array}$$


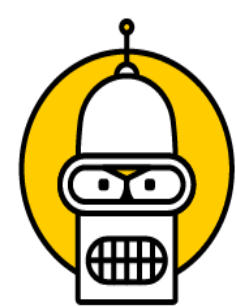
The requester and performers maximize their utility simultaneously if the pricing $p(a)$ for each label is proportional to its accuracy a .

Performance-based pricing in practice: settings

- › Price p for the level of accuracy a_0 : $\Pr(\hat{z} = z) \geq a_0$ E.g.:



- › $\hat{q}_w = \Pr(y^w = z)$ - estimated quality level of performer w ,
e.g. the fraction of correct labels for golden set (GS):



5 correct GS
among 10
 $\hat{q}_w = 0.5$



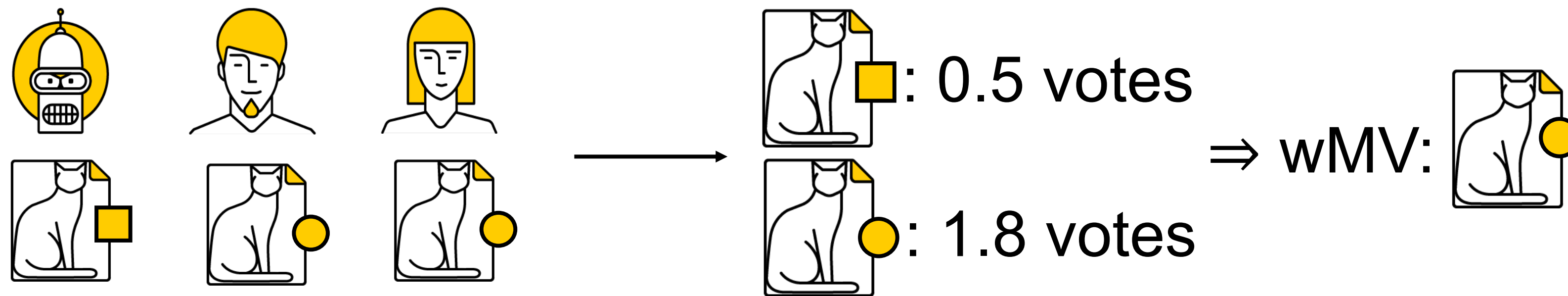
16 correct GS
among 20
 $\hat{q}_w = 0.8$



100 correct GS
among 100
 $\hat{q}_w = 1$

Performance-based pricing in practice: settings

› Aggregation $\hat{z}_j^{\text{wMV}} = \arg \max_{y=1,\dots,K} \sum_{w \in W_j} \hat{q}_w \delta(y = y_j^w)$



› IRL algorithm is based on the expected accuracy of \hat{z}_j^{wMV}

Performance-based pricing in practice

Pricing rules

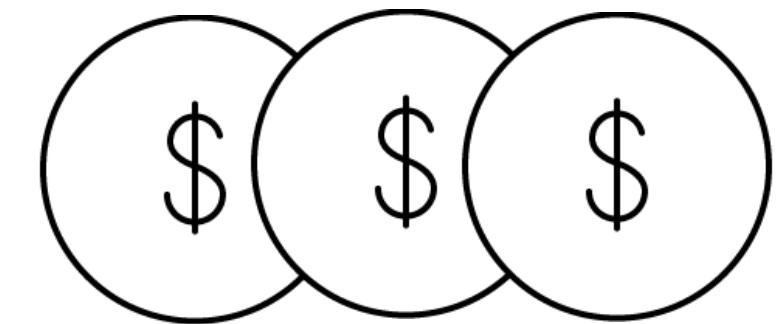
1. If $\hat{q}_w \geq a_0$, then the price is p
2. Else find n :

$$\underbrace{\sum_{k=0}^{n/2} \binom{n}{k} \hat{q}_w^{n-k} (1 - \hat{q}_w)^k}_{\text{Expected accuracy for MV}} \geq a_0$$

Expected accuracy for MV

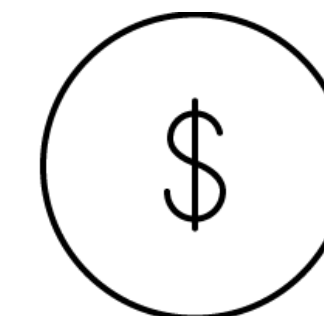
The price is p/n

$$a_0 = 0.99$$



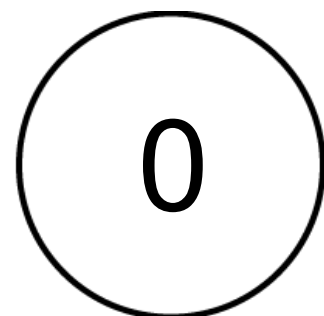
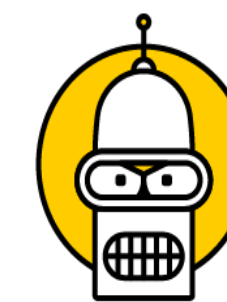
0.3\$

$$\hat{q}_w = 1$$



0.02\$

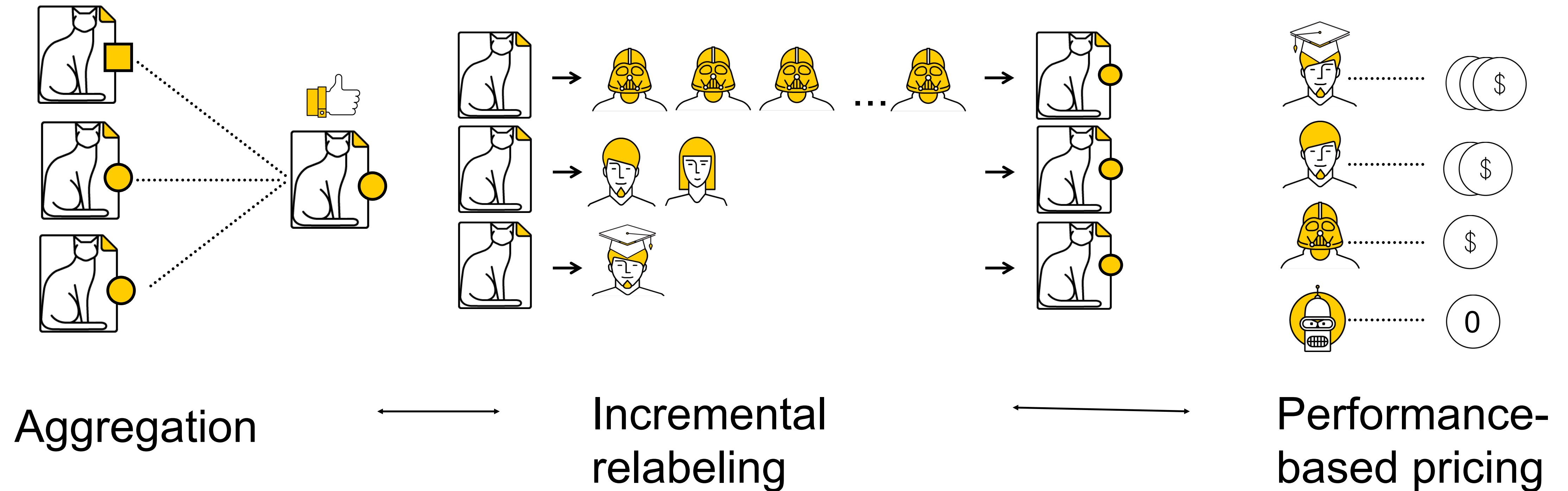
$$\hat{q}_w = 0.8$$
$$\Rightarrow n = 15$$



0\$

$$\hat{q}_w = 0.5$$
$$\Rightarrow n = \infty$$

Key components of labeling with crowds



Tutorial Schedule

```
graph TD; A[Introduction: 15 min] --> B[Part I: 30 min  
Key Components for Data Collection]; B --> C[Part II: 60 min  
Practice Session I]; C --> D[Lunch Break: 45 min]; D --> E[Part III: 45 min  
Advanced Techniques]; E --> F[Part IV: 30 min  
Practice Session II]; F --> G[Part V: 15 min  
Conclusion];
```

Introduction: 15 min

Part I: 30 min
Key Components for
Data Collection

Part II: 60 min
Practice Session I

Lunch Break:
45 min

Part III: 45 min
Advanced
Techniques

Part IV: 30 min
Practice Session II

Part V: 15 min
Conclusion



Thank you!
Questions?

Dmitry Ustalov

Analyst/Software Developer
Crowdsourcing Research Group



dustalov@yandex-team.ru



<https://research.yandex.com/tutorials/crowd/naacl-2021>