

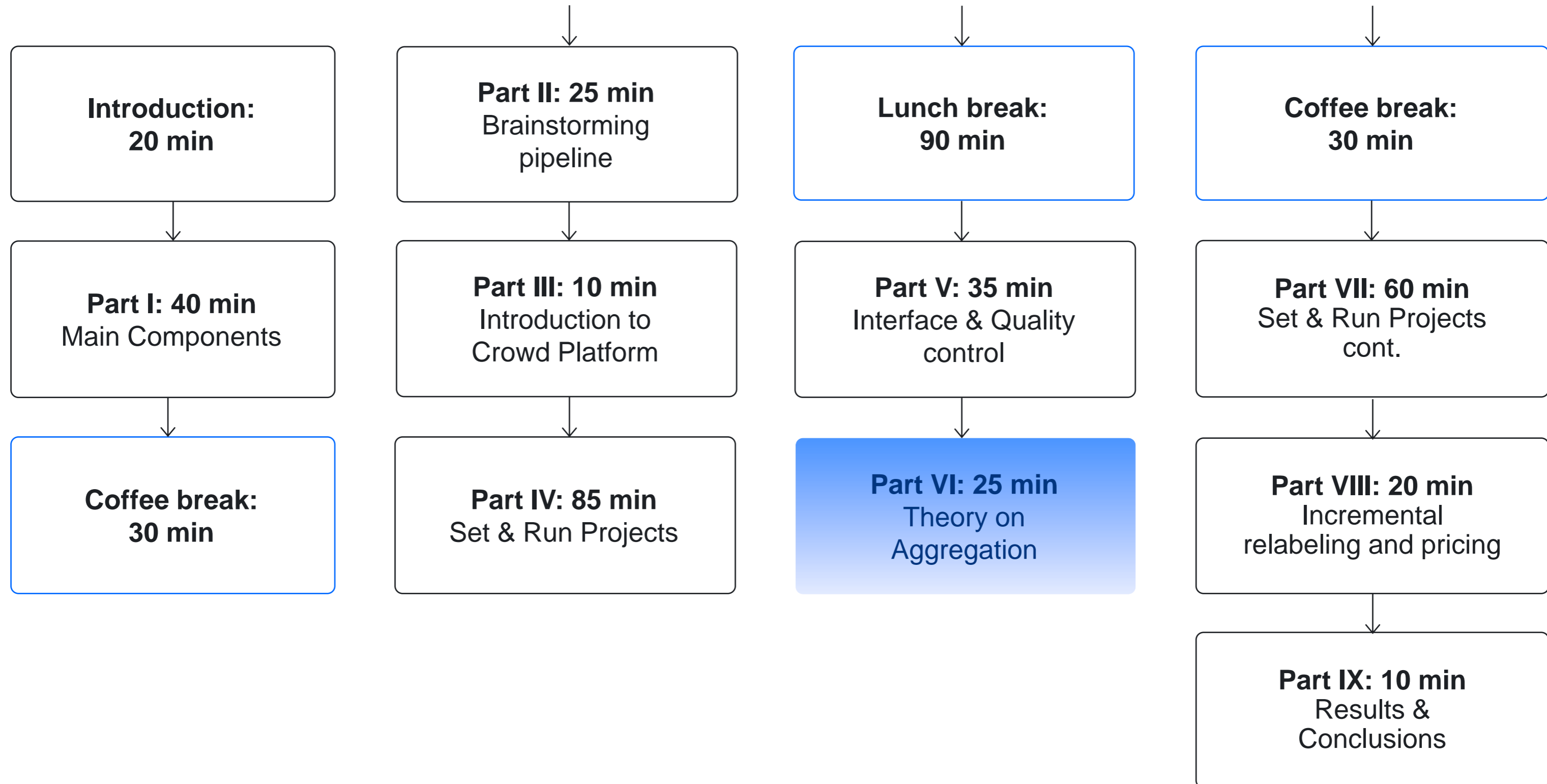
Part VI

Theory on Aggregation

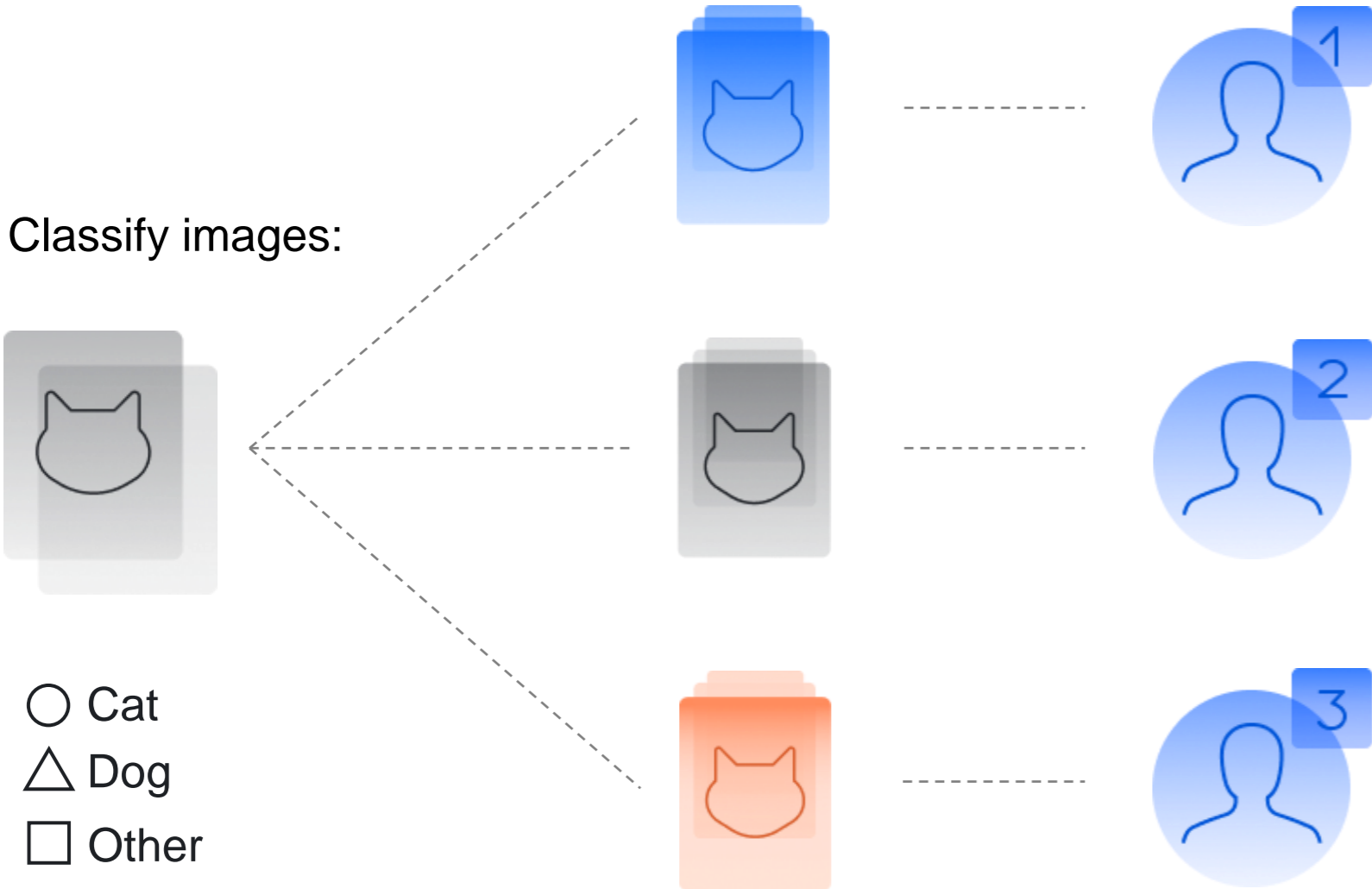
Valentina Fedorova,
Research analyst

Toloka

Tutorial schedule



Labeling data with crowdsourcing



- ▶ How to choose a reliable label?
- ▶ How many workers per object?
- ▶ How much to pay to workers?
- ▶ ...

Evaluation of labeling approaches



VS

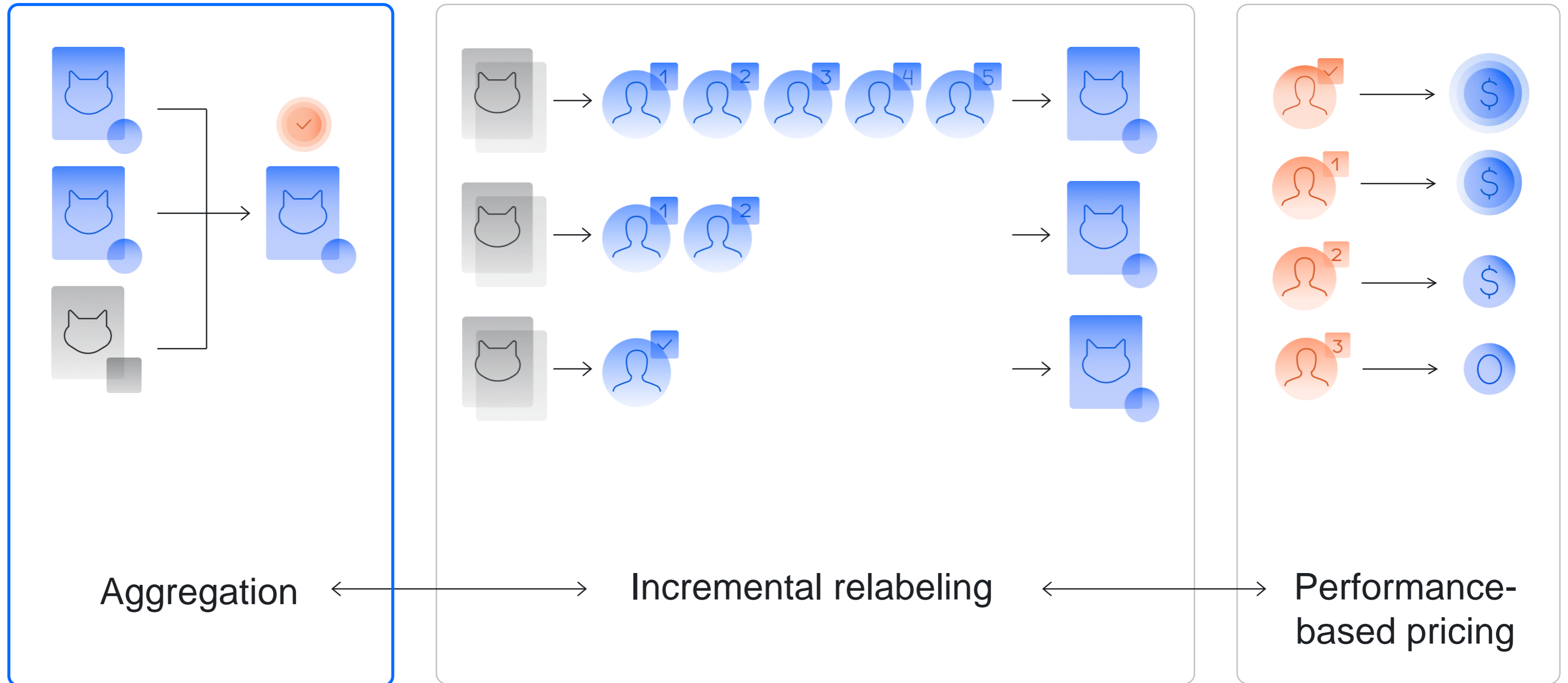


Accuracy

Cost

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

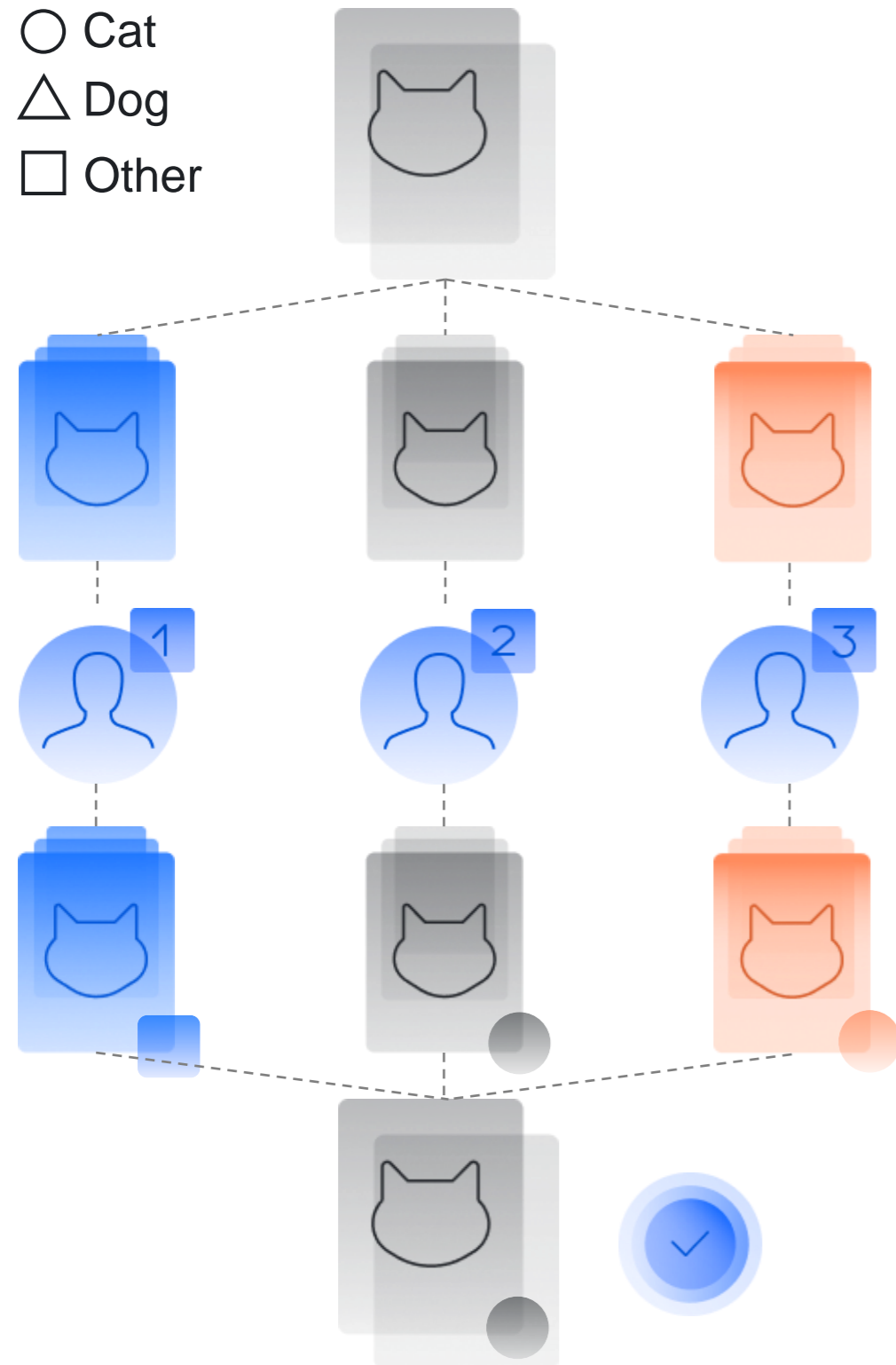
Key components of labeling with crowds



Aggregation

The background features a series of overlapping, curved, blue shapes that create a sense of depth and movement, resembling a stylized sunburst or a series of concentric arcs. The colors range from a deep navy blue to a bright, vibrant blue.


Labeling data with crowds





- ▶ Classify images
- ▶ Upload multiple copies of each object to label
- ▶ Workers assign noisy labels to objects
- ▶ Aggregate multiple labels for each object into a more reliable one

Process results




Projects > Does the image contains traffic lights? > pool

 **pool** — closed ▲

Statistics Download results  ^ Edit  ?

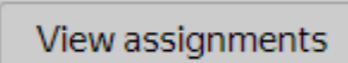
- View operations
- Dawid-Skene aggregation model
- Aggregation by skill


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



0  30

Multiclass labels

Project 1: Filter images

Are there shoes
in the picture?

Yes

No



Notation

► Categories $k \in \{1, \dots, K\}$. E.g.:

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

► Workers: $w \in \{1, \dots, W\}$. E.g.:

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

○ Cat

△ Dog

□ Other



The simplest aggregation: Majority Vote (MV)

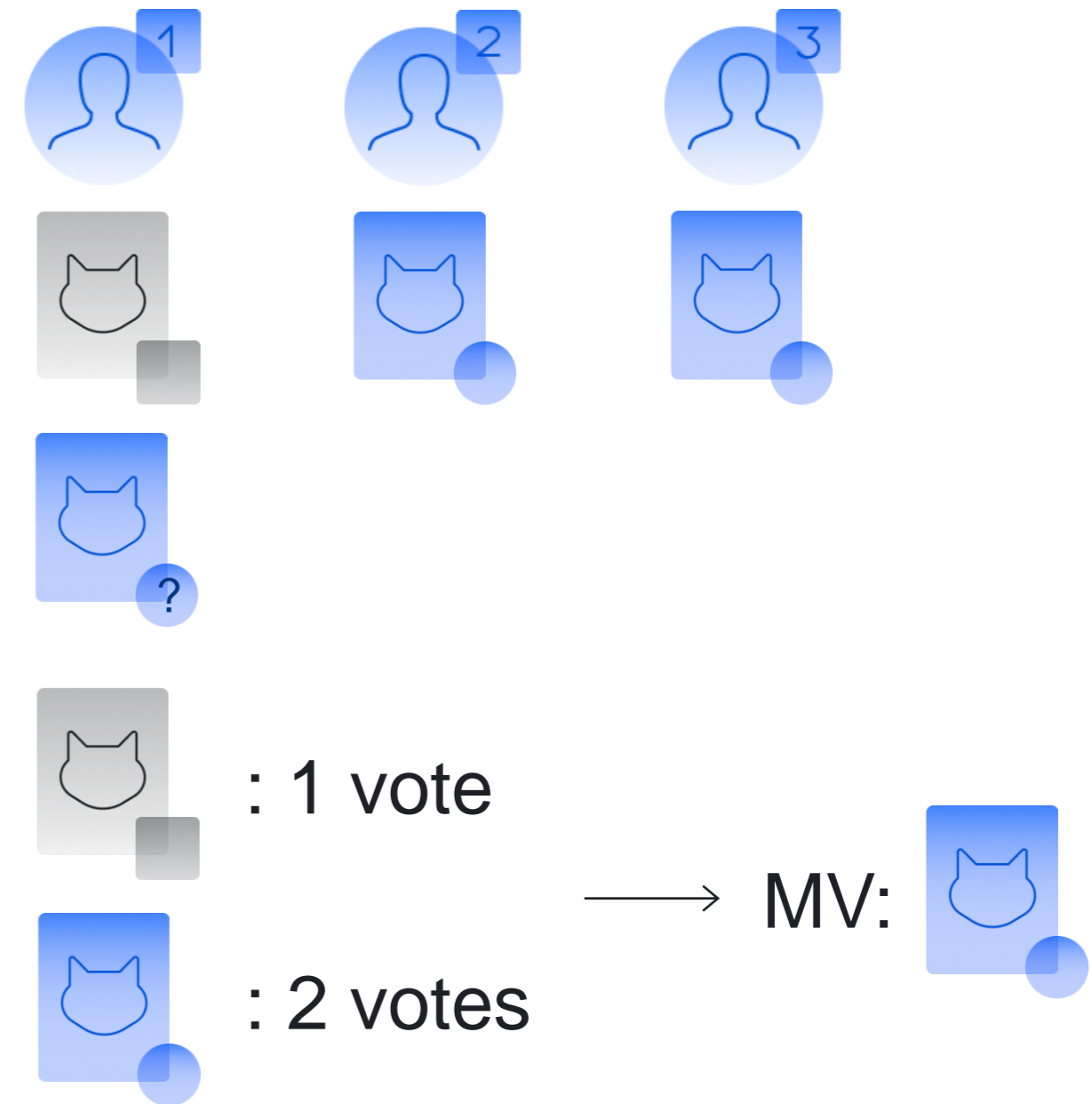
► The problem of aggregation:

- Observe noisy labels

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

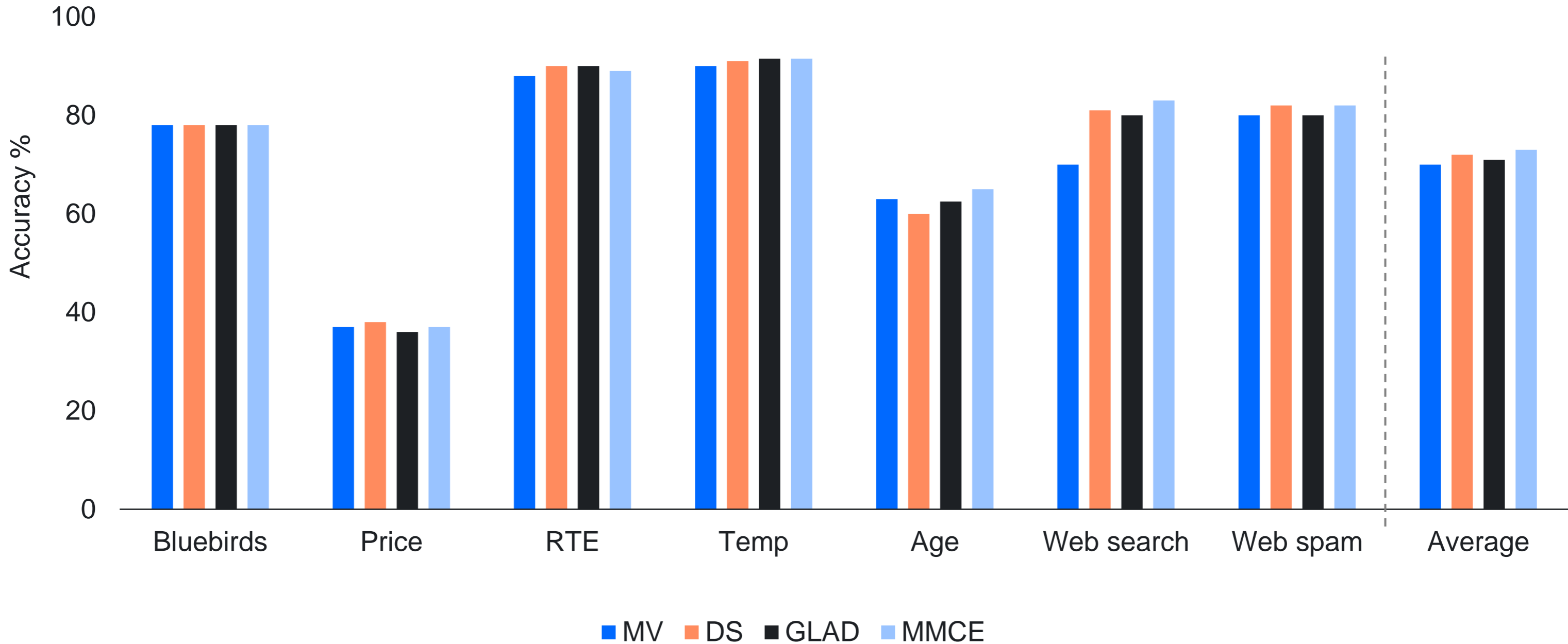
- Recover true labels $z = \{z_j \mid j = 1, \dots, J\}$

► 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

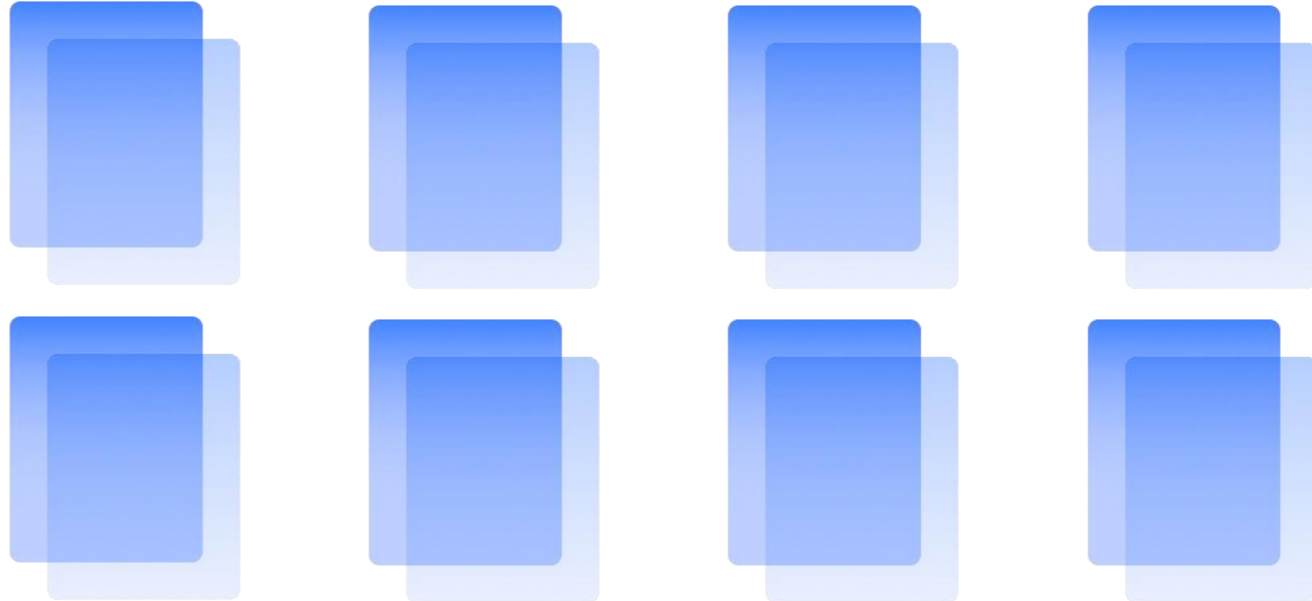


Properties of MV

All workers are treated similarly

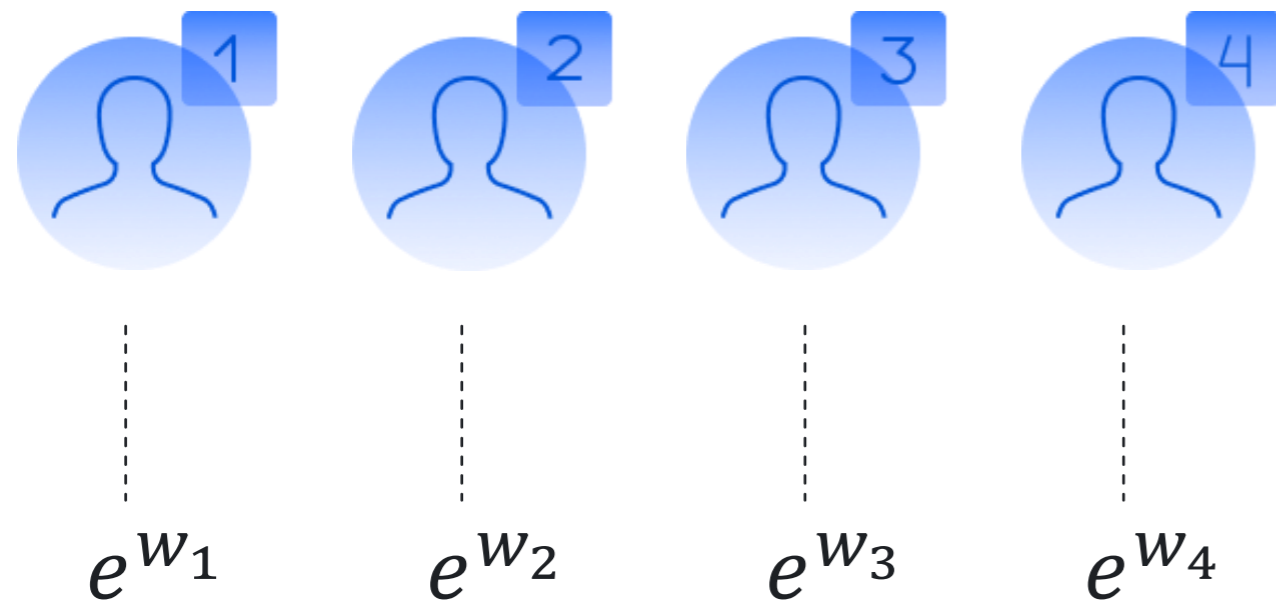


All objects are treated similarly

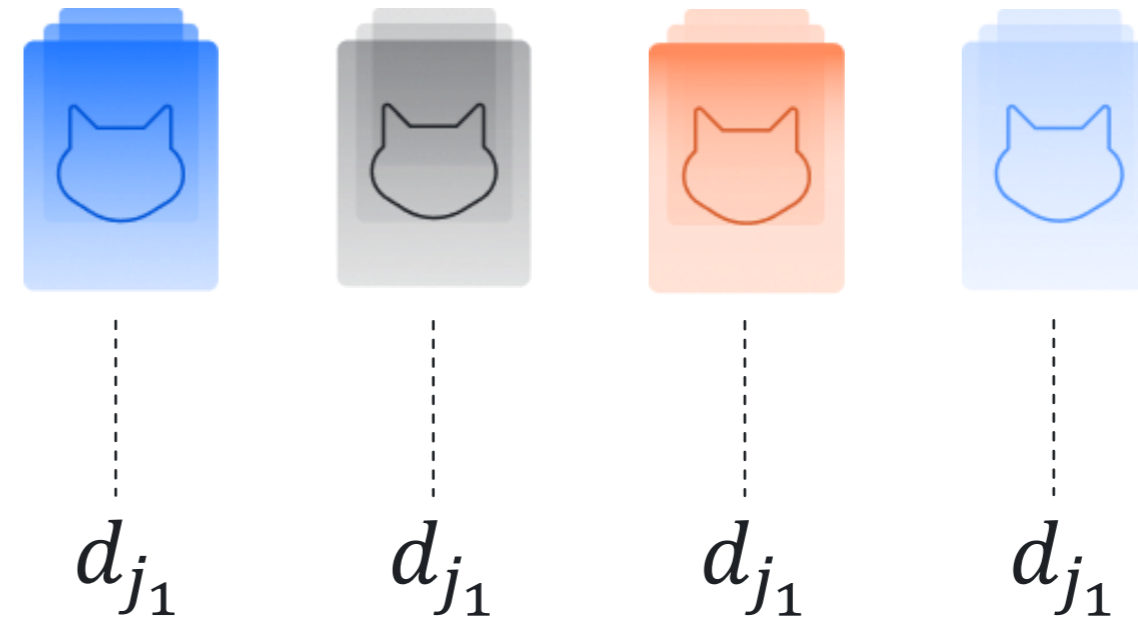


Advanced aggregation: workers and objects

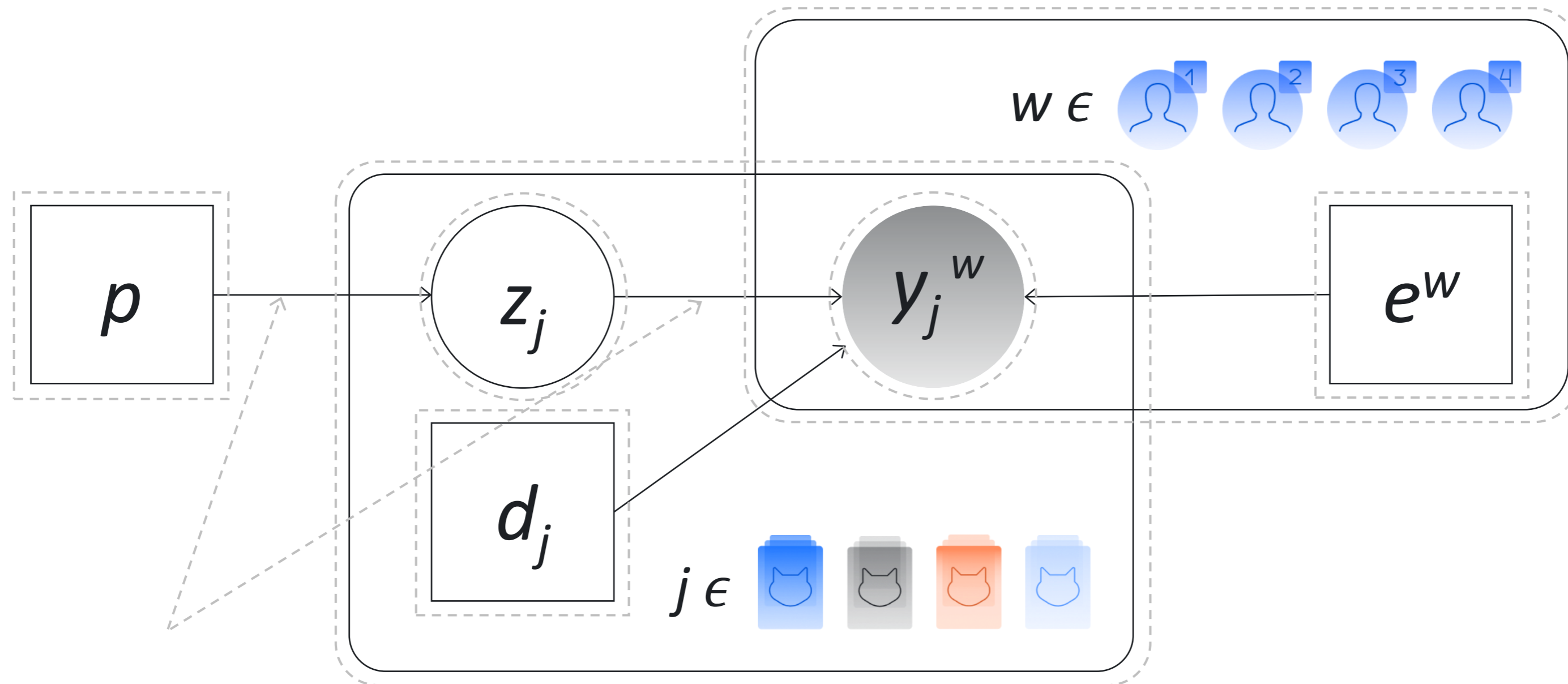
Parameterize expertise of workers by e^w



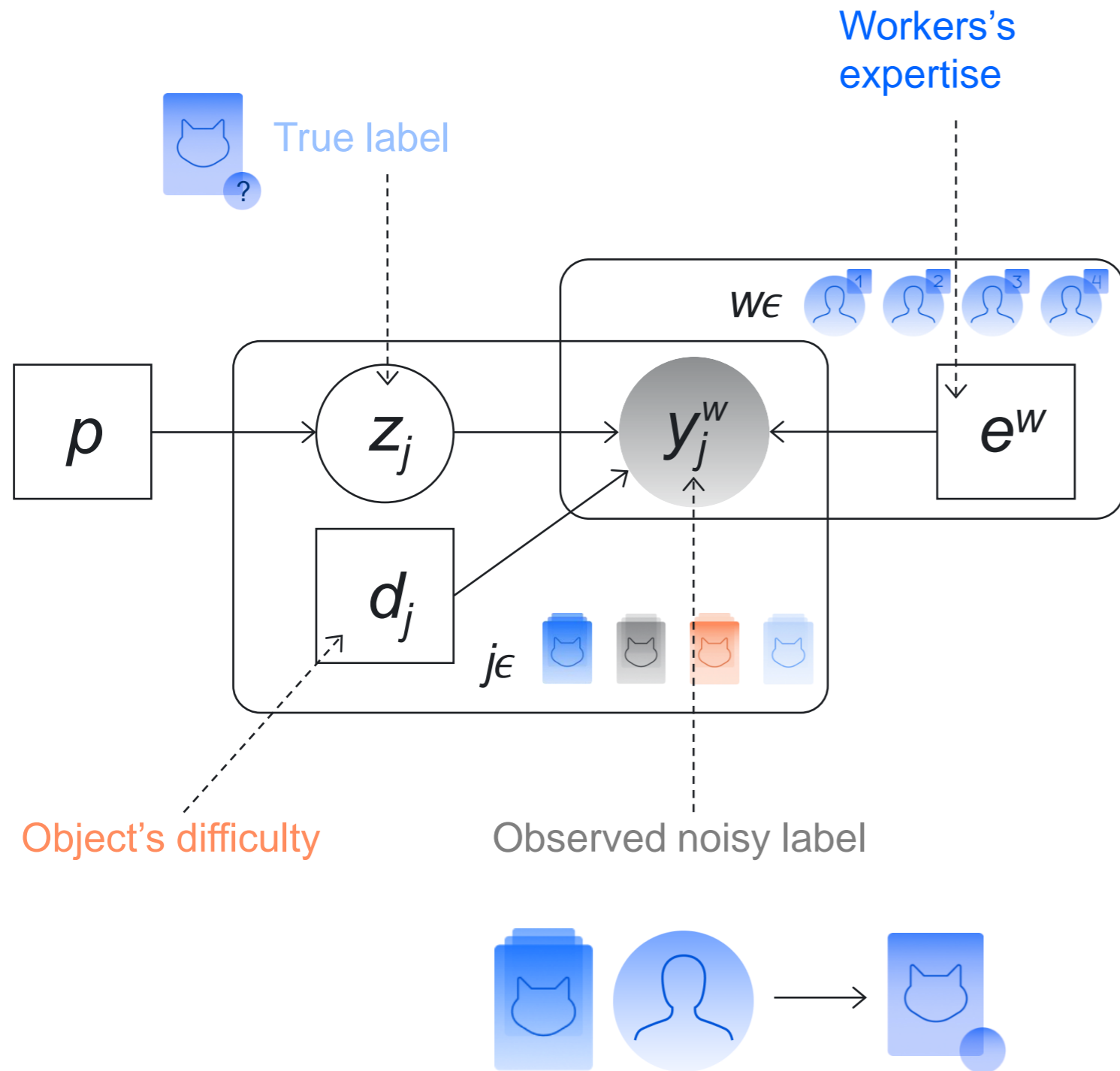
Parameterize difficulty of objects by d_j



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

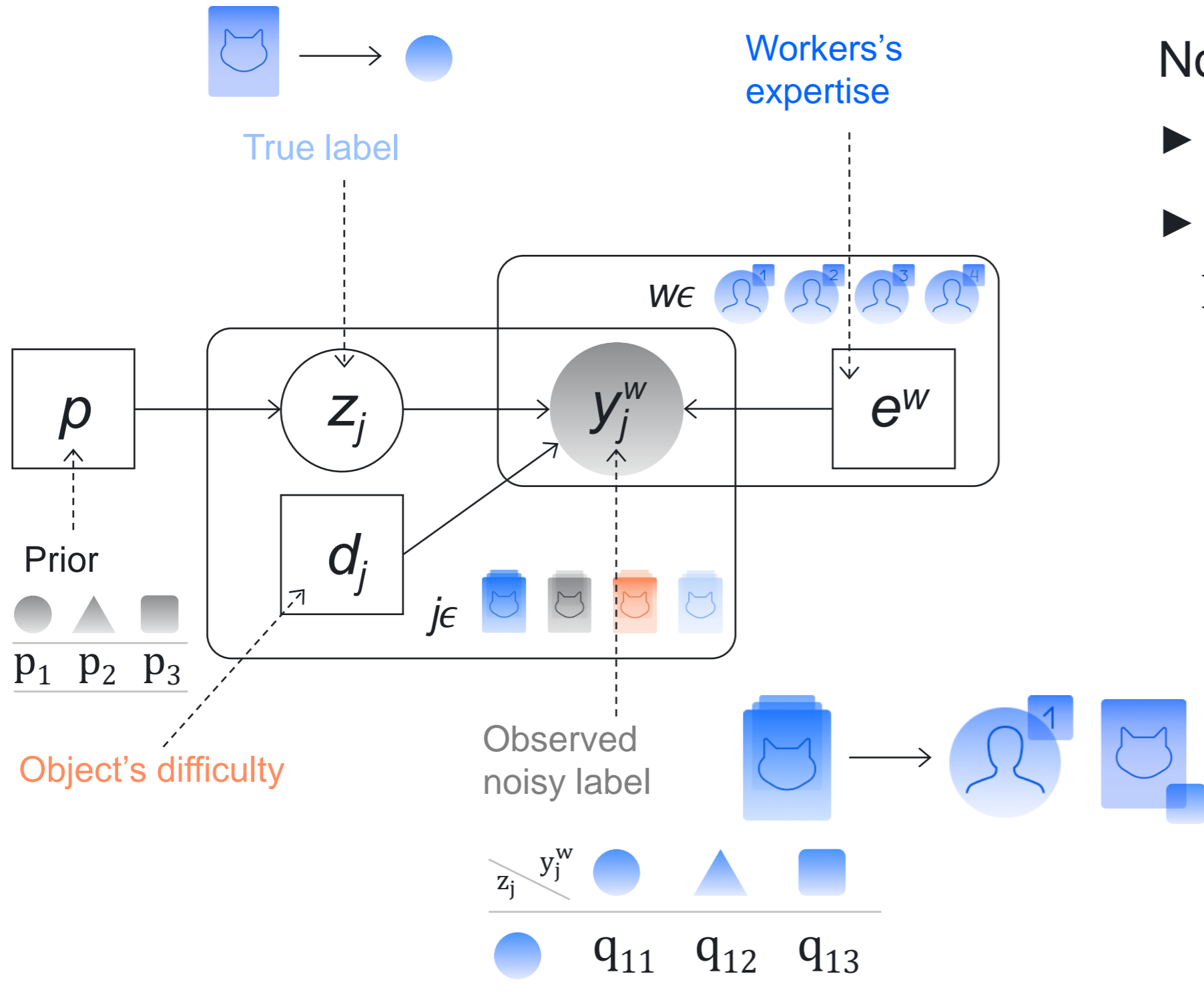
$$M_j^w[c, k] = \Pr(Y_j^w = k | Z_j = c)$$



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:

- ▶ Sample z_j from a distribution $P_Z(p)$
- ▶ 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 $\text{Mult}(\cdot)$

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 y and z under the latent label model:

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

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)*

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

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

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

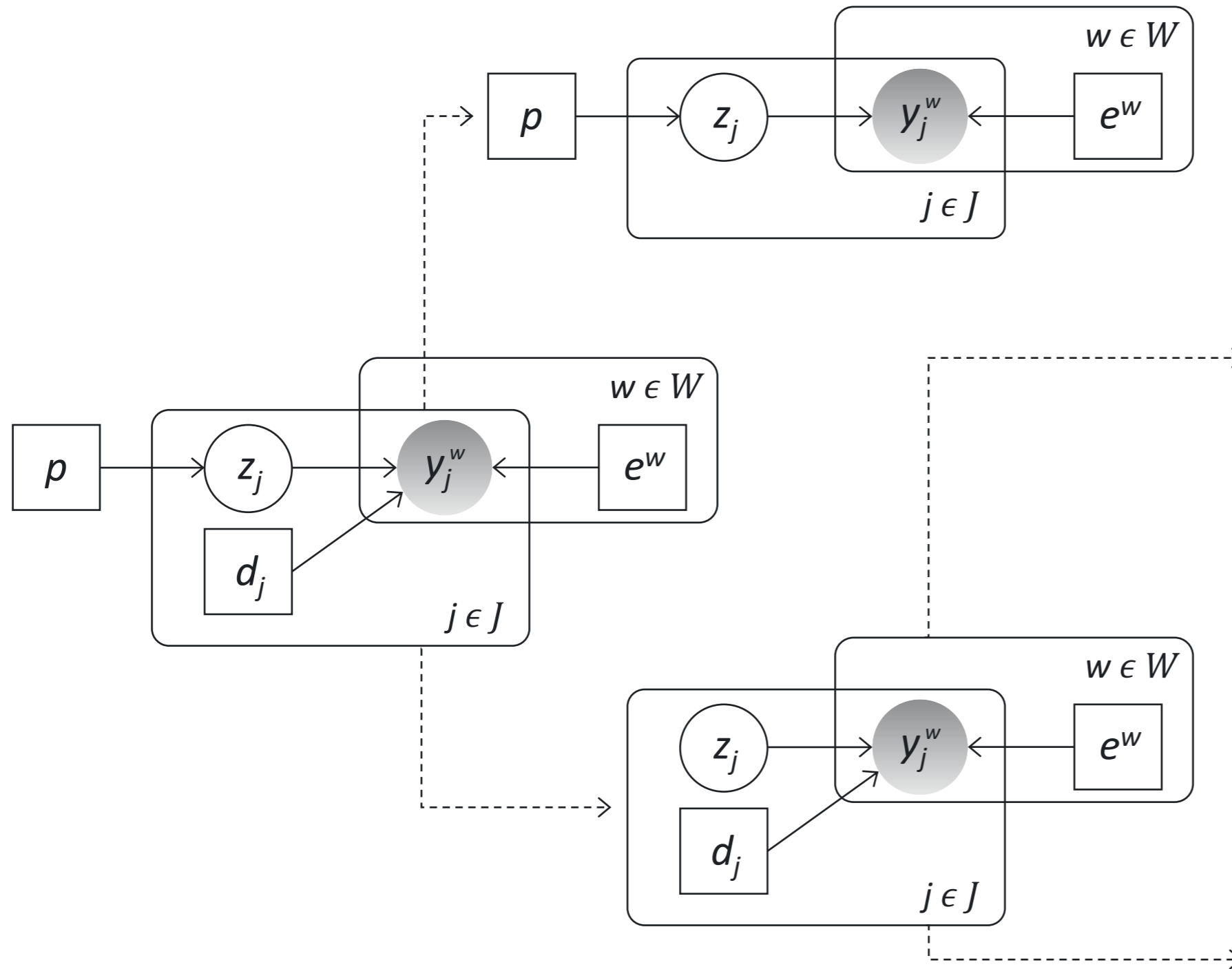
- ▶ **M-step:** Maximize the expectation of LL with respect to the posterior distribution of \hat{z} :

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

- Analytical solutions
- Gradient descent

* it is a lower bound on LL of \mathbf{y} and \mathbf{z}

Latent label model (LLM): special cases

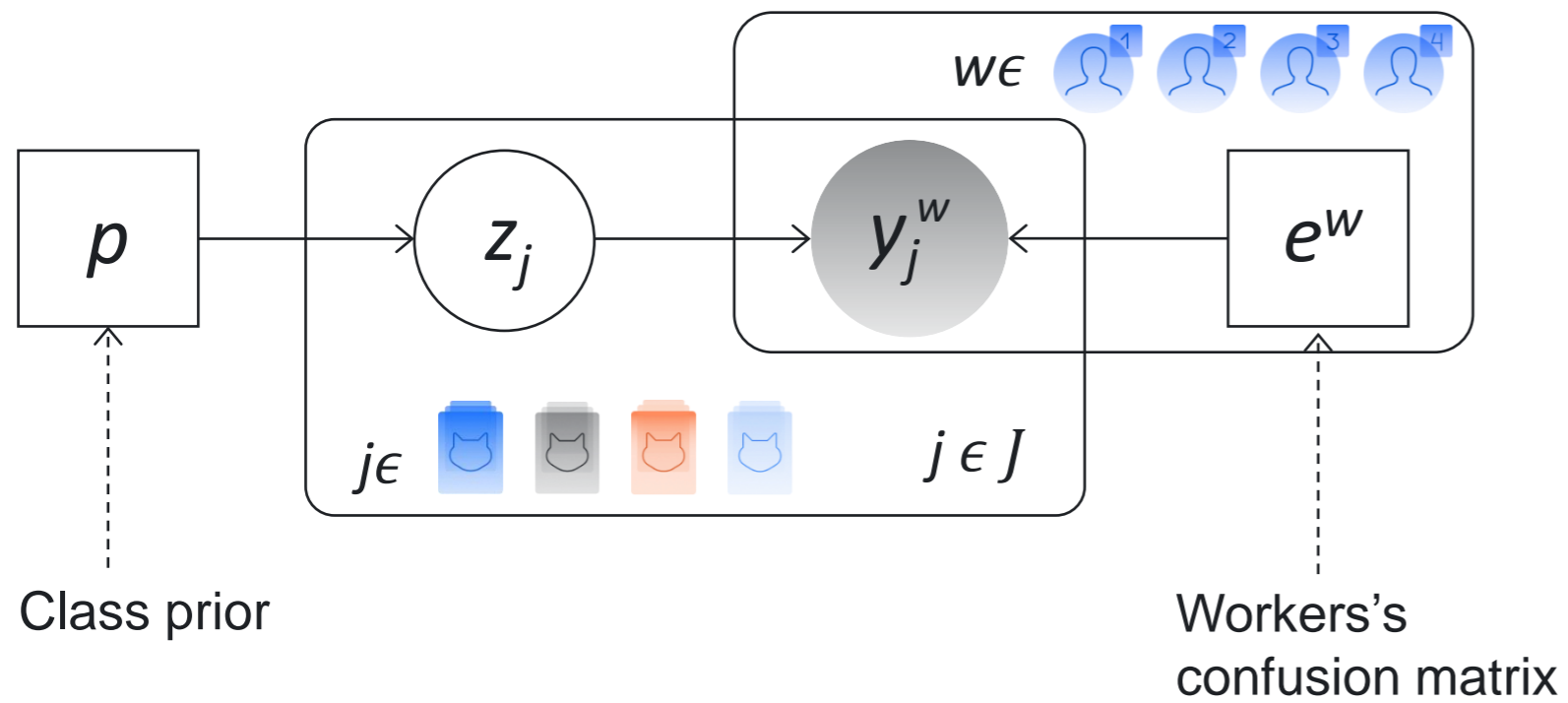


- ▶ Dawid and Skene model (DS):
 - Categories are **different**
 - Objects are **similar**
 - Workers are **different**

- ▶ Generative model of labels, abilities, and difficulties (GLAD):
 - Categories are **similar**
 - Objects are **different**
 - Workers are **different**

- ▶ Minimax conditional entropy model (MMCE):
 - Categories are **different**
 - Objects are **different**
 - Workers are **different**

Dawid and Skene model (DS)



LLM with parameters:

- ▶ p — vector of length K : $p[i] = \Pr(Z = c)$
- ▶ e^w — matrix of size $K \times K$: $e^w[c, k] = \Pr(Y^w = k | Z = c)$

$z \backslash y_w$	●	▲	■
●	■	■	■
▲	■	■	■
■	■	■	■

DS: parameters optimization

► **E-step:**

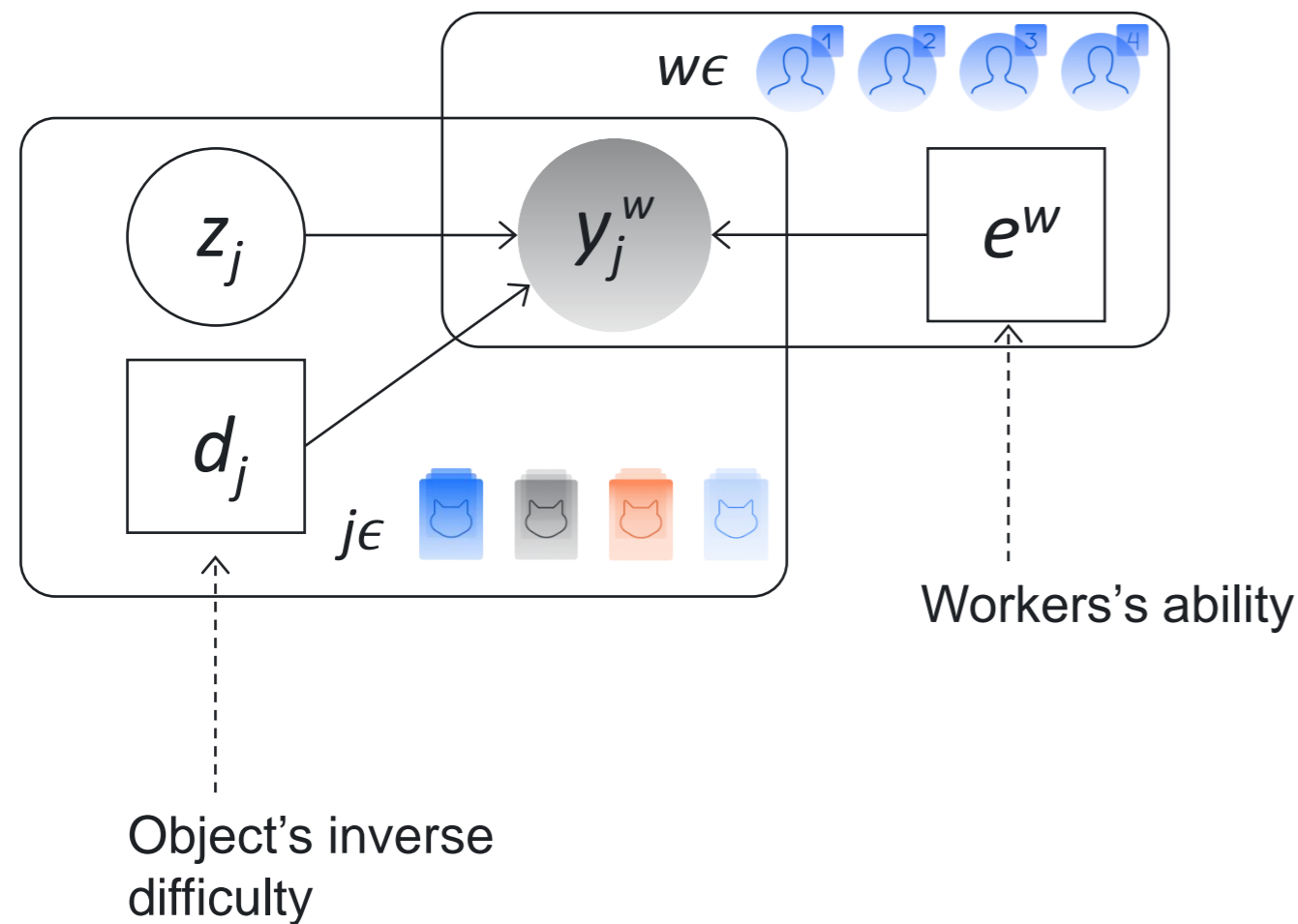
$$\hat{z}_j[c] = \frac{p[c] \prod_{w \in W_j} e^w[c, y_j^w]}{\sum_k p[k] \prod_{w \in W_j} e^w[k, y_j^w]}, \quad c = 1, \dots, K$$

► **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., $P(z_j) = 1/K$)

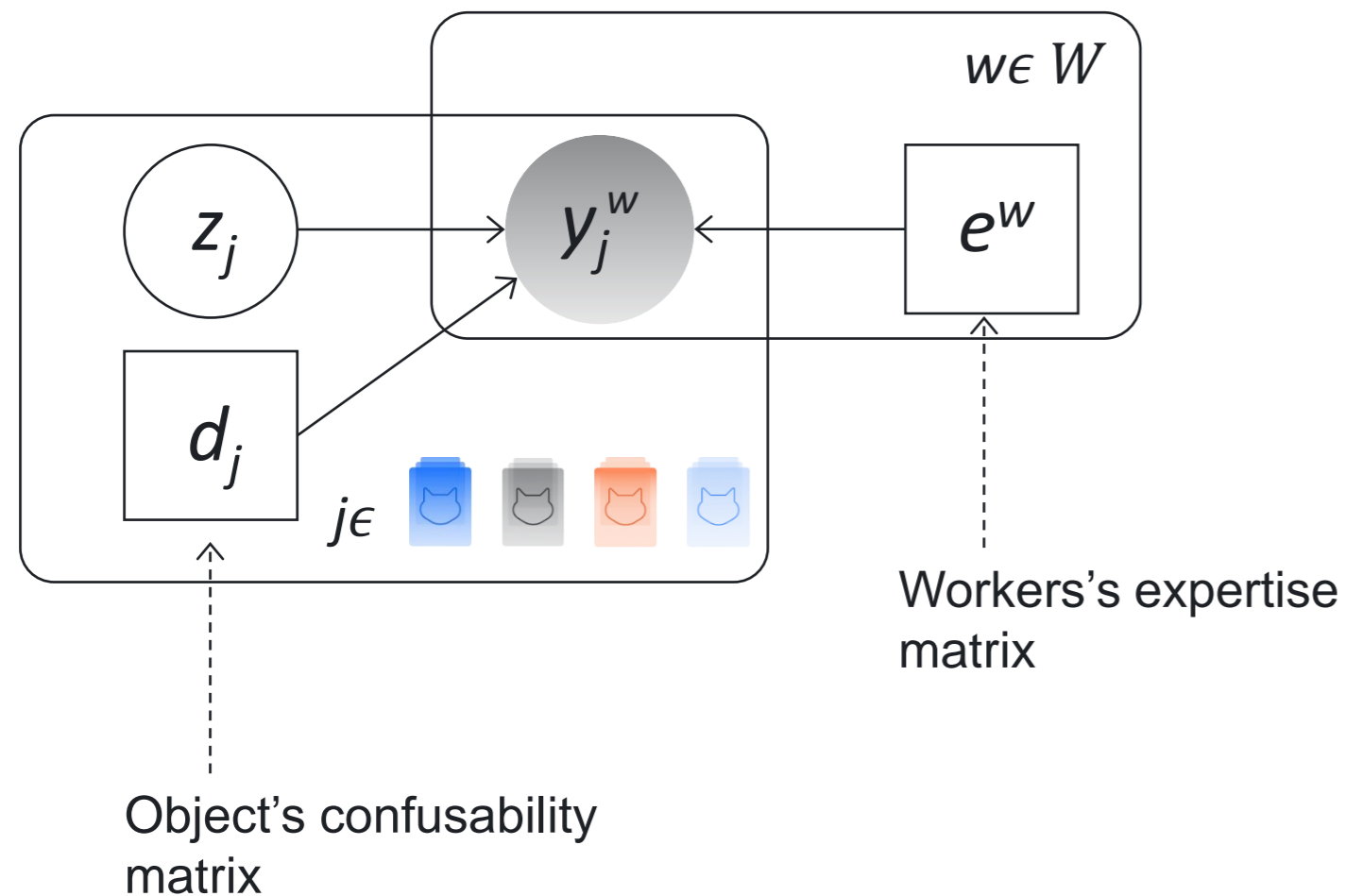
► **E-step:**

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

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

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

MiniMax Conditional Entropy model (MMCE)







































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

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

- 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)												
Workers (W)												
Number of parameters	0			$WK^2 + K$			$W + J$			$(W + J)K^2$		

Pairwise comparisons

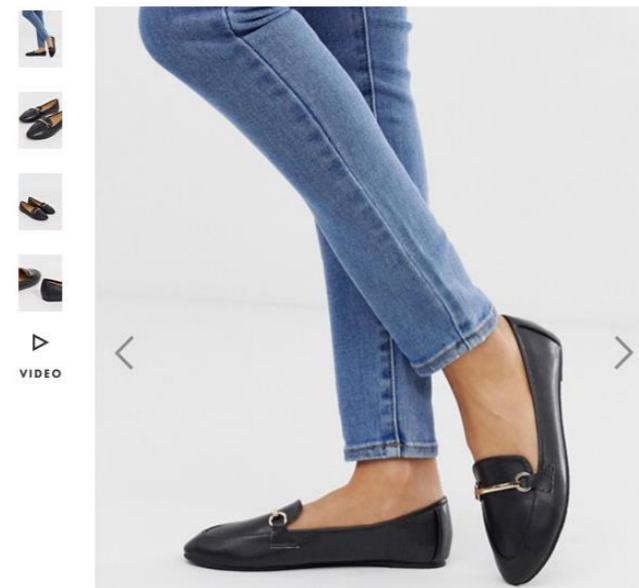
The background of the slide is a dark blue gradient with a series of concentric, curved lines that create a sense of depth and movement, resembling a stylized eye or a series of overlapping planes.

Project 4: Compare items

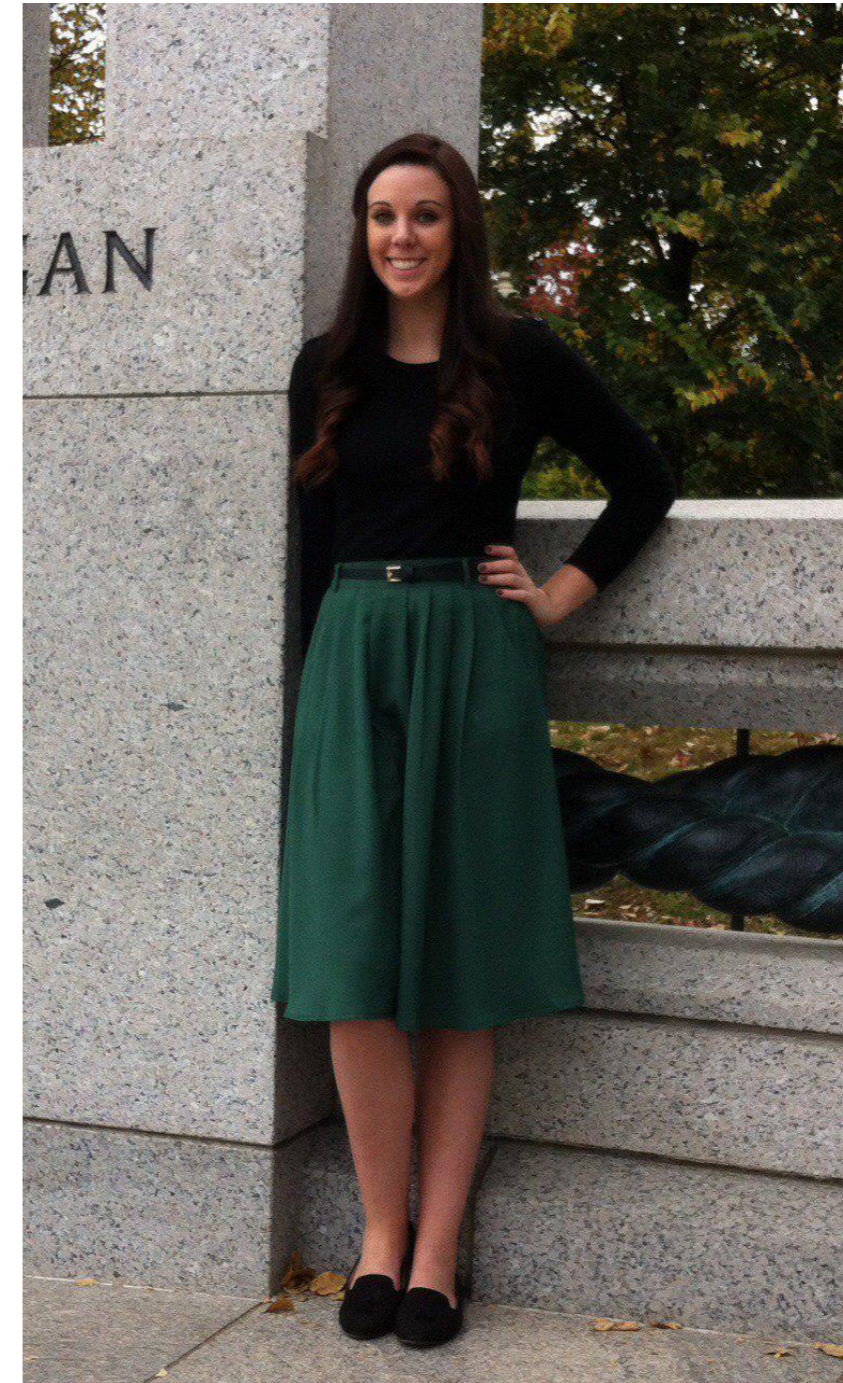
Which shoes look more similar to the one in the picture?



Left



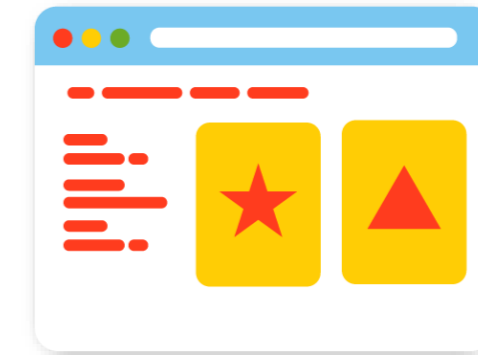
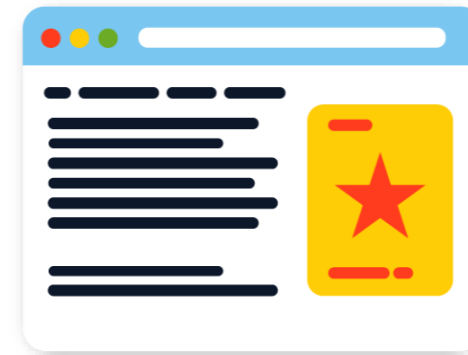
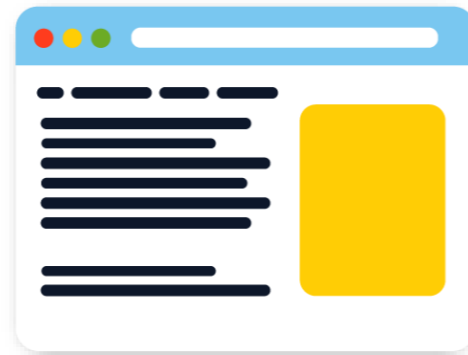
Right



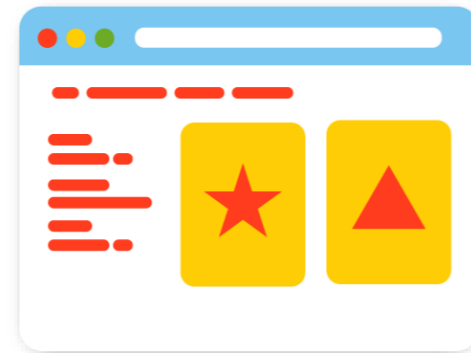
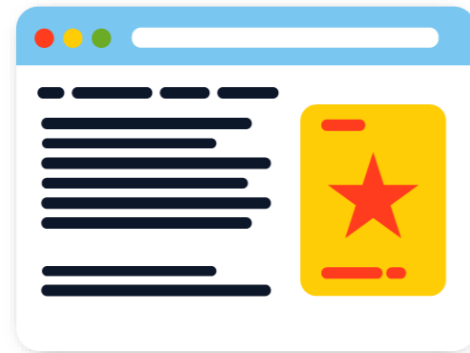
Notation

▶ Answers: **Left** or **Right**

▶ Items $d_j \in \{1, \dots, N\}$ E.g.:



▶ Tasks:



Choose a better item:

Left
Right

▶ Workers $w \in \{1, \dots, W\}$ E.g.:



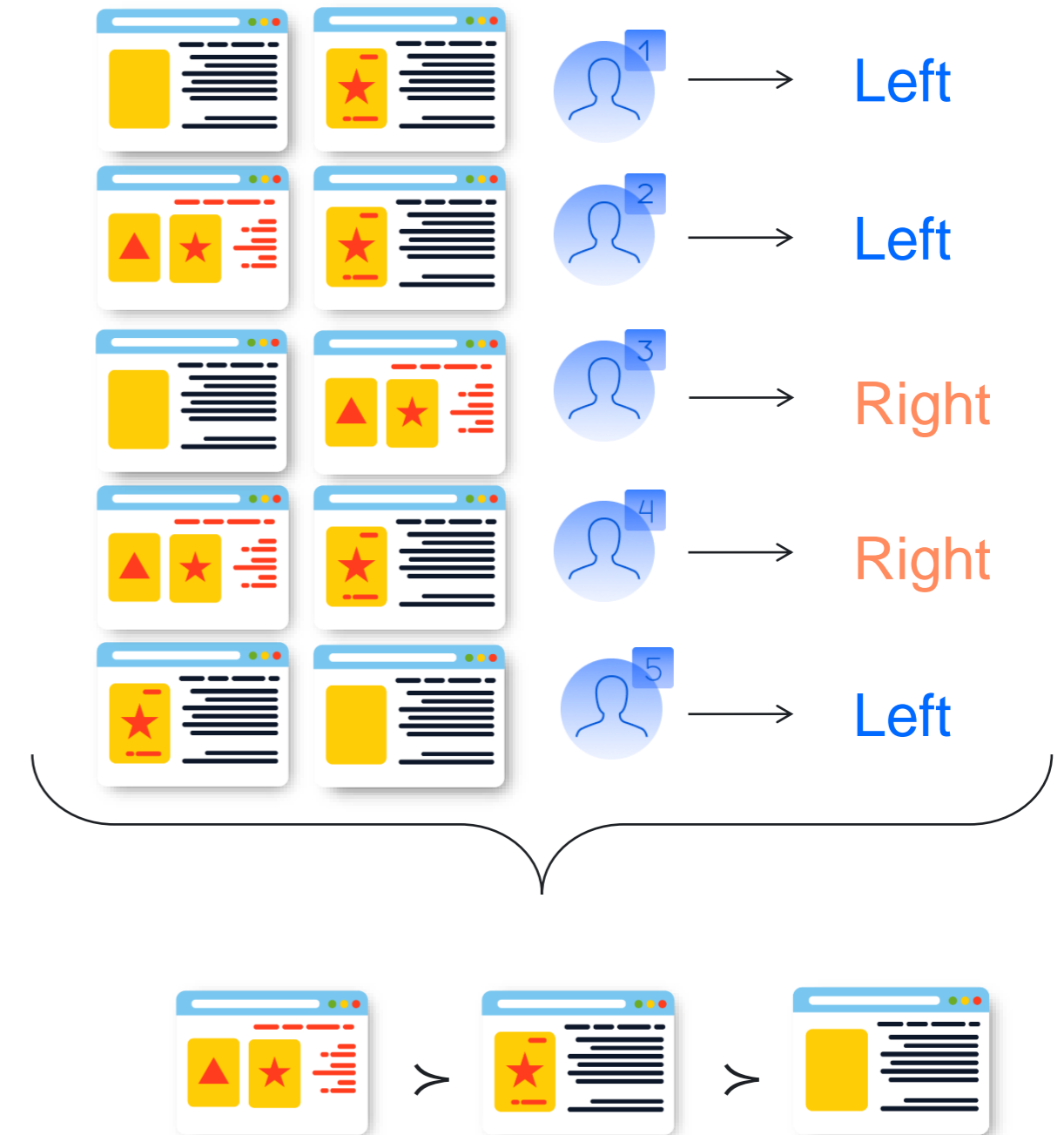
Formalization

Ranking from pairwise comparisons:

- ▶ Given pairwise comparisons for items in D :

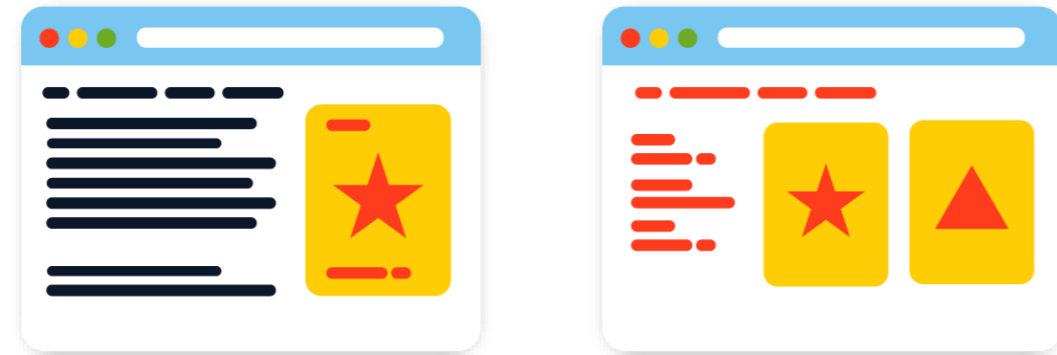
$$P = \{(w_k, d_i, d_j) : i \succ_k j\}$$

- ▶ Obtain a **ranking** π over items $D \rightarrow \{1, \dots, N\}$ based on answers in P



Difference from multiclassification

- ▶ The latent label assumption is not satisfied when comparing complex items

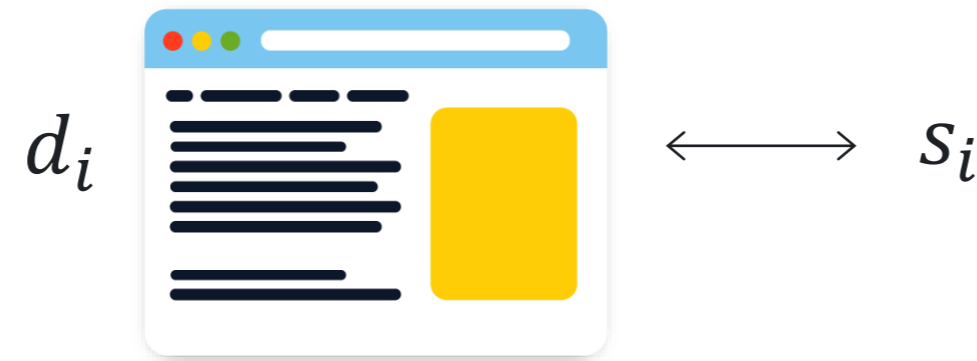


- ▶ Different tasks may contain common items



Bradley and Terry model (BT)

- ▶ Assume that each item $d_i \in D$ has a latent “quality” score $s_i \in \mathbb{R}$



- ▶ The probability that $d_i \in D$ will be preferred in a comparison over $d_j \in D$

$$\Pr(i \succ j) = f(s_i - s_j), \text{ where } f(x) = 1/(1+e^{-x}).$$

- ▶ The model assumes that all workers are equally good and truthful

NoisyBT model: parameterization of workers

w_k  \longleftrightarrow “reliability” γ_k and “bias” q_k

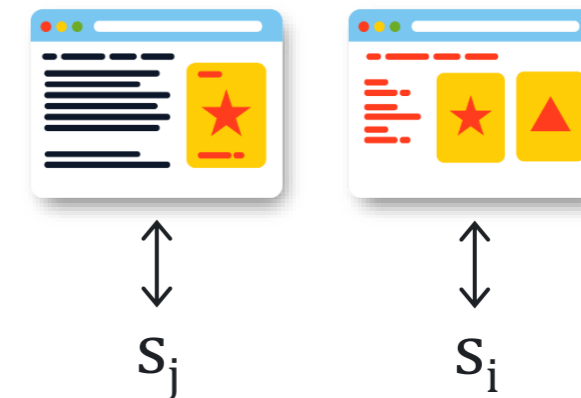
- ▶ The probability that w reads task is

$$\Pr(w_k \text{ reads a task}) = f(\gamma_k) \leftarrow \text{Logistic function}$$

- ▶ If w_k reads a task, she answers according to scores:

$$(f(s_i - s_j), f(s_j - s_i))$$

Probability to choose **Left** if compares items



- ▶ If w_k does not read a task, she answers according to her bias

$$(f(q_k), f(-q_k))$$

Probability to choose **Left** if answers randomly

NoisyBT: likelihood of workers' answers

The likelihood of $i \succ_k j$ is

$$\Pr(i \succ_k j) = \underbrace{f(\gamma_k) f(s_i - s_j)}_{\text{Truthful answer}} + \underbrace{\left(1 - f(\gamma_k)\right) f\left((-1)^{(1 - \mathbb{I}(d_i \text{ was left}))} q_k\right)}_{\text{Random answer}},$$

where $\mathbb{I}(d_i \text{ was left})$ is the indicator for the order of d_i and d_j



NoisyBT: parameters optimization

Likelihood of observed comparisons:

$$T(s, q, \gamma) = \sum_{(w_k, d_i, d_j) \in P} \log \Pr(i \succ_k j) =$$

$$\sum_{(w_k, d_i, d_j) \in P} \log [f(\gamma_k) f(s_i - s_j) + (1 - f(\gamma_k)) f((-1)^{(1 - \mathbb{I}(d_i \text{ was left}))} q_k)]$$

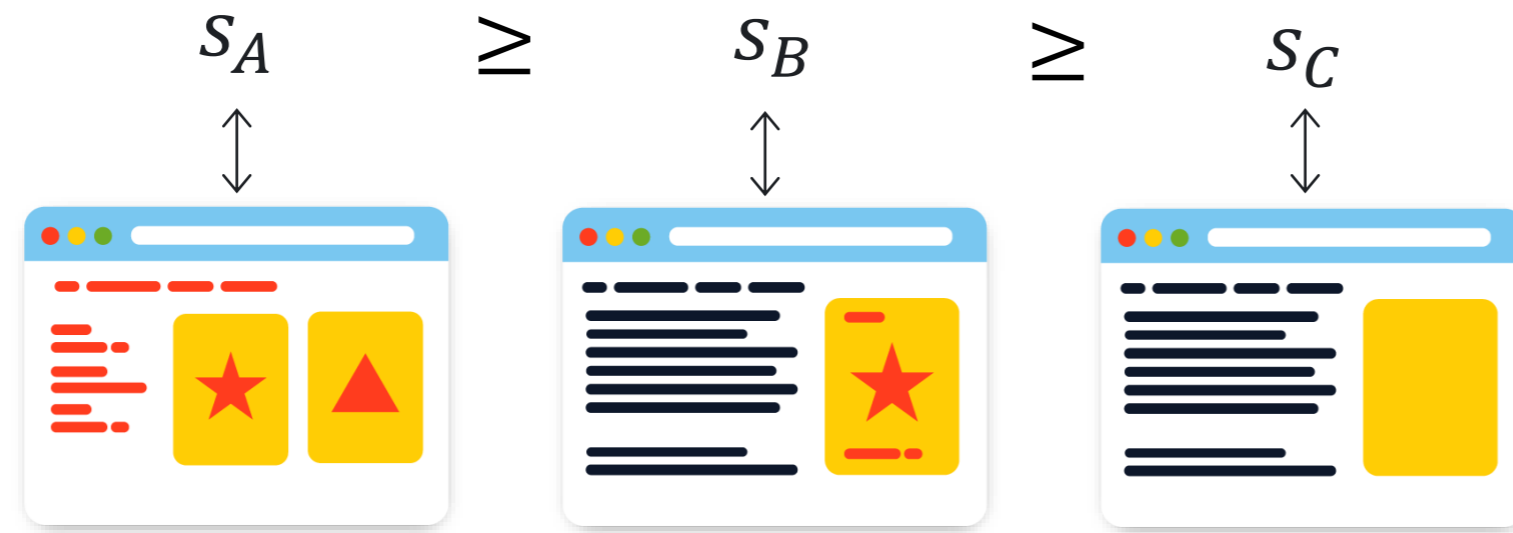
- ▶ $\{s_i\}_{i=1, \dots, N}$ and $\{\gamma_k, q_k\}_{k=1, \dots, W}$ are inferred by maximizing the log-likelihood:

$$T(s, q, \gamma) \rightarrow \max_{\{s_i, \gamma_k, q_k\}}$$

- ▶ To obtain a **ranking** π over items, **sort** items according to their **scores**

Summary about pairwise comparisons

- ▶ Latent scores models for ranking from pairwise comparisons:



- ▶ To reduce bias from unreliable answers parameterize workers

