



Toloka

# Improving Web Ranking with Human-in-the-Loop: Methodology, Scalability, Evaluation

Alexey Drutsa, Dmitry Ustalov, Nikita Popov,  
Daria Baidakova



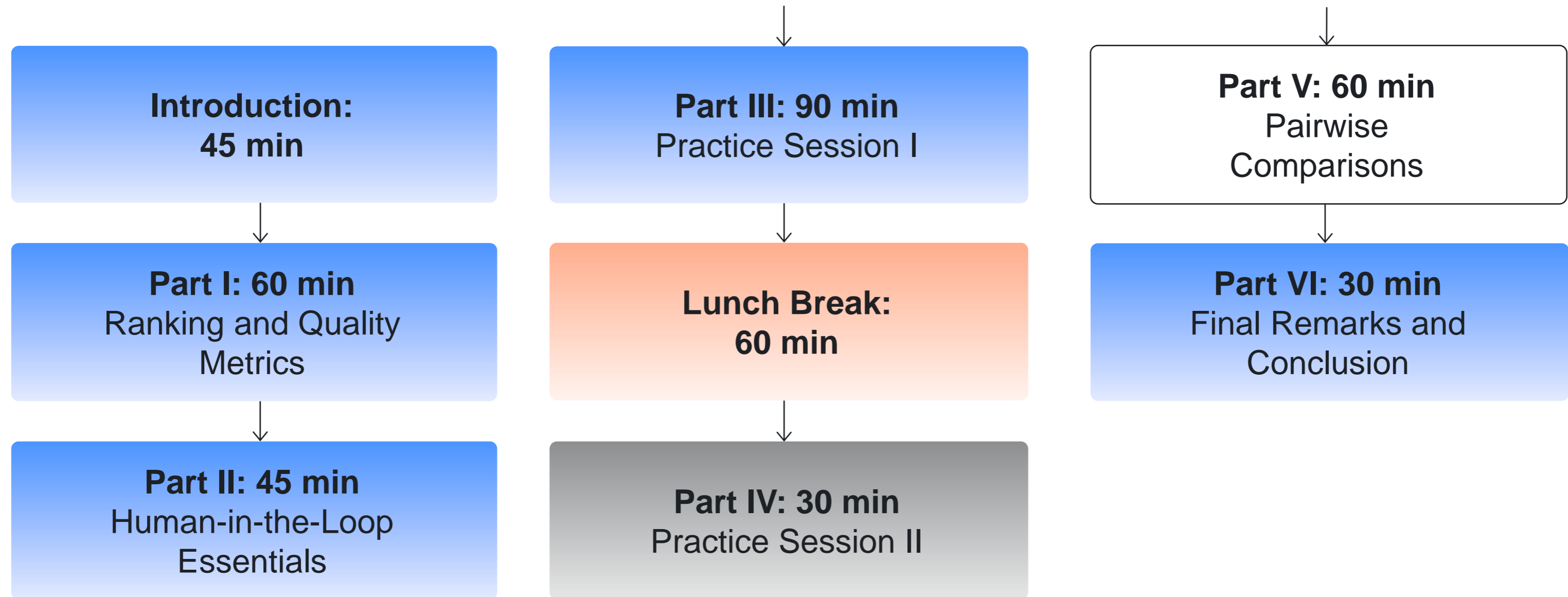
# Part V

# Pairwise

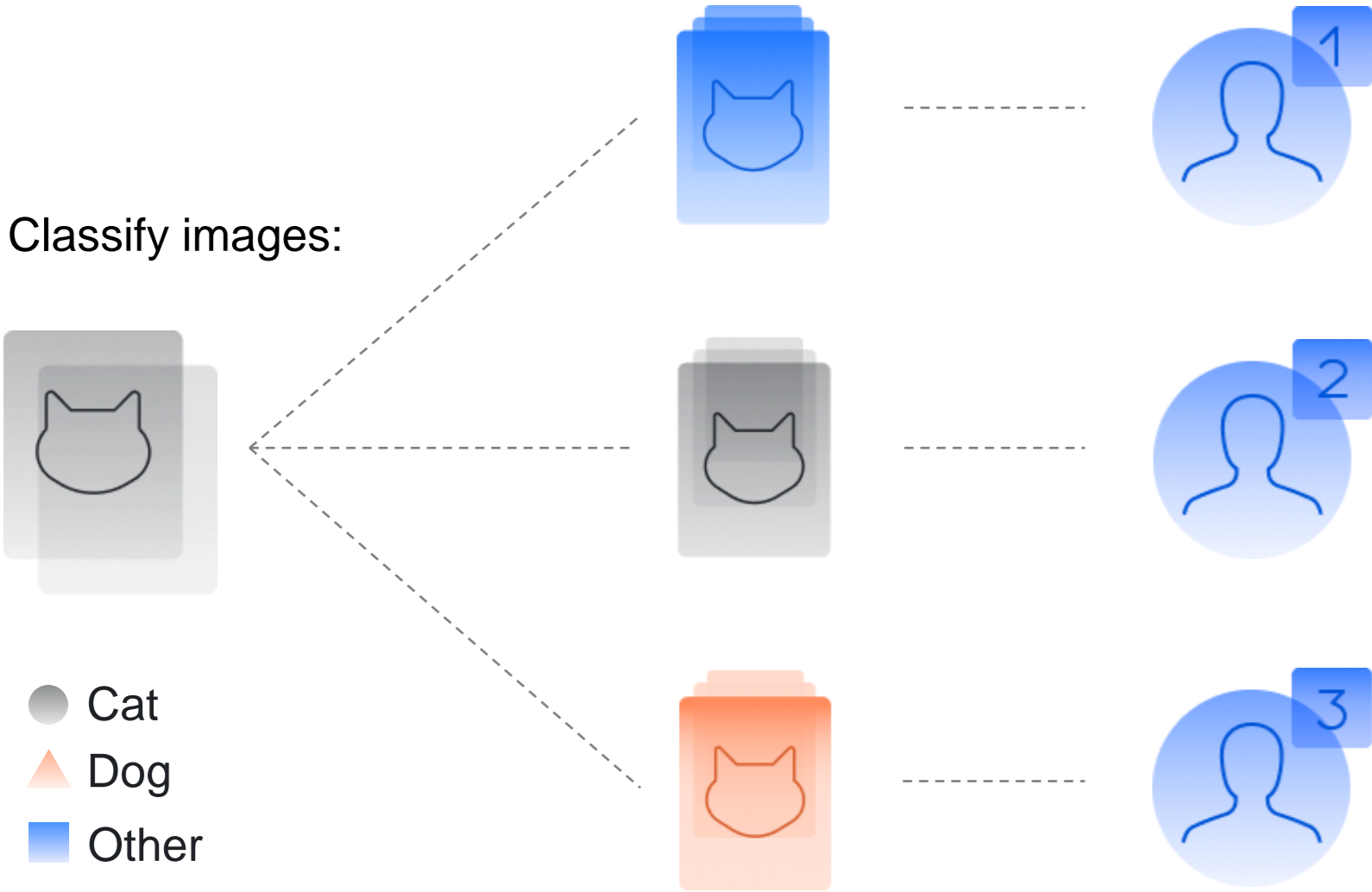
# Comparisons

Dmitry Ustalov,  
Analyst/Software Developer

# Tutorial schedule



# Labelling Data with Crowdsourcing



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

# Evaluation of Labelling Approaches



VS



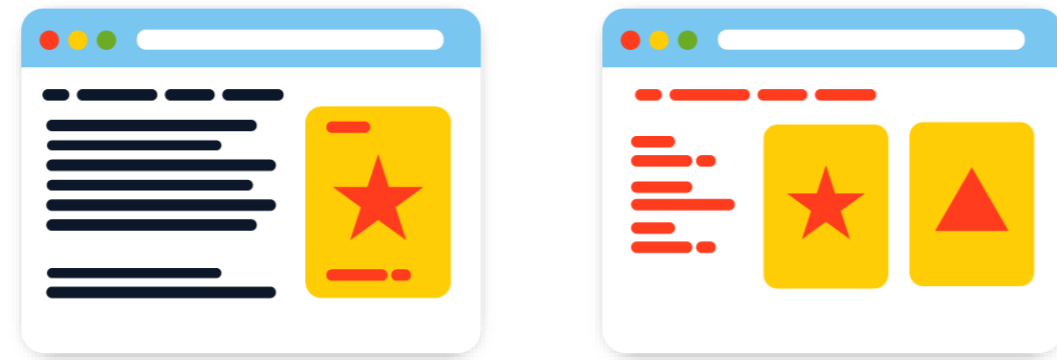
Accuracy

Cost

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

# Difference from Multiclassification

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



- ▶ Different tasks may contain common items

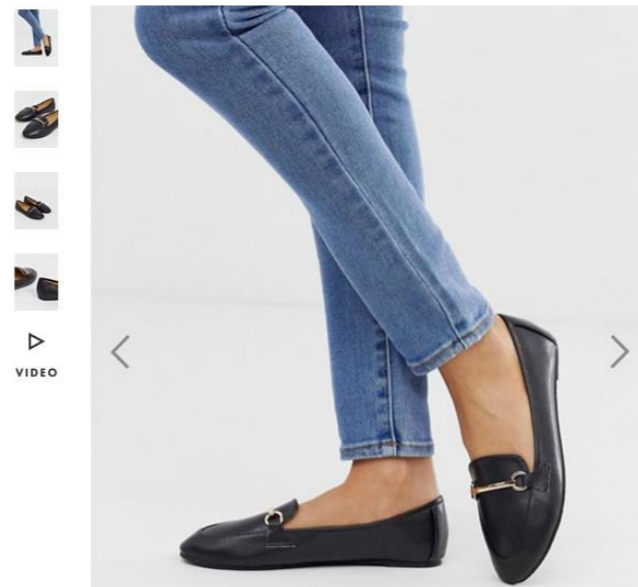


# Task: Compare Items

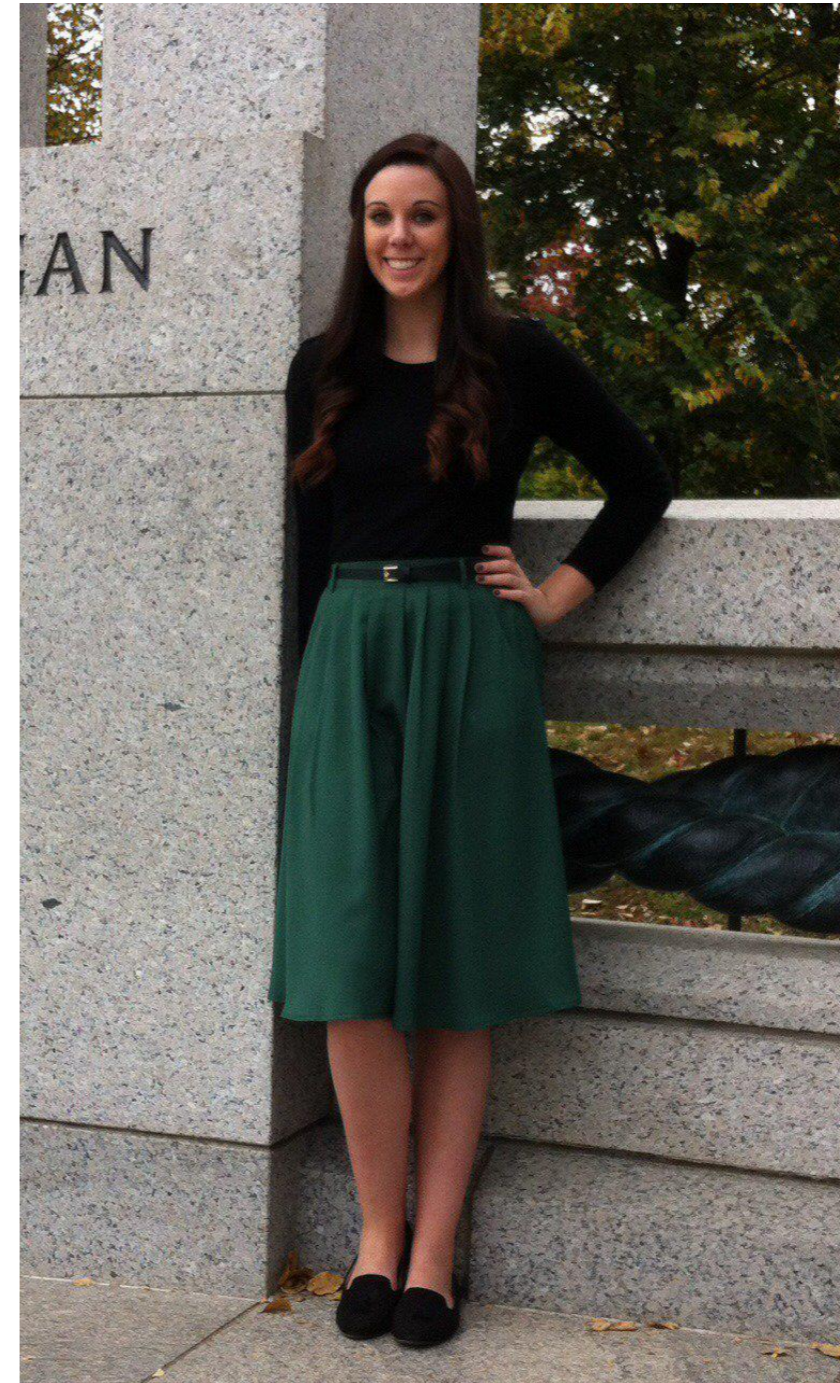
Which shoes look more similar to the one in the picture?



Left



Right



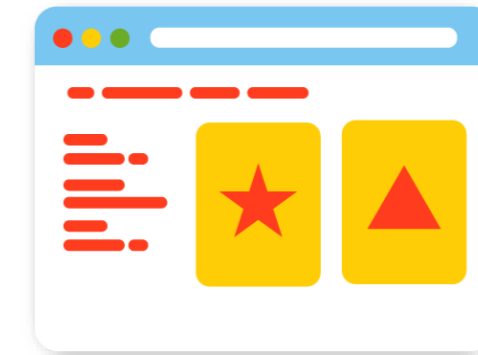
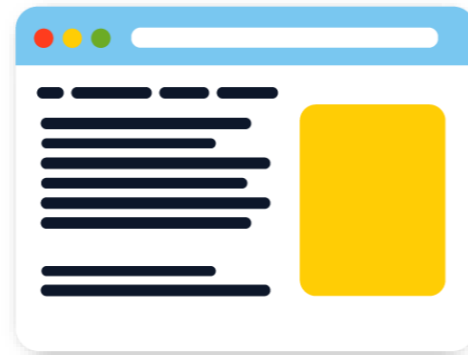


# Aggregating Pairwise Comparisons

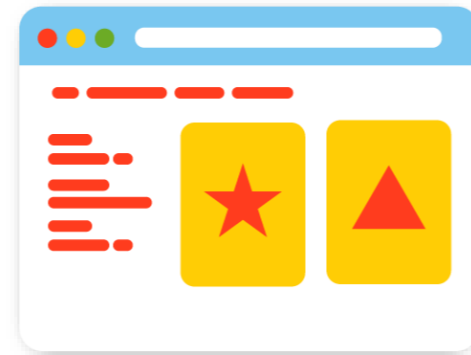
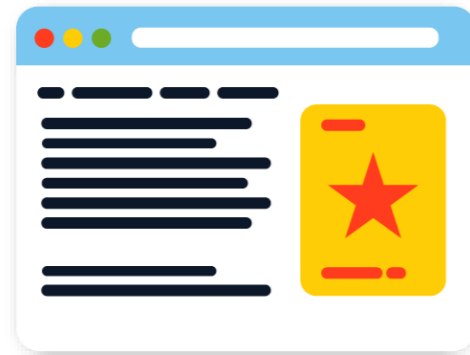
# Notation

▶ Answers: **Left** or **Right**

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



▶ Tasks:



Choose a better item:

**Left**  
**Right**

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



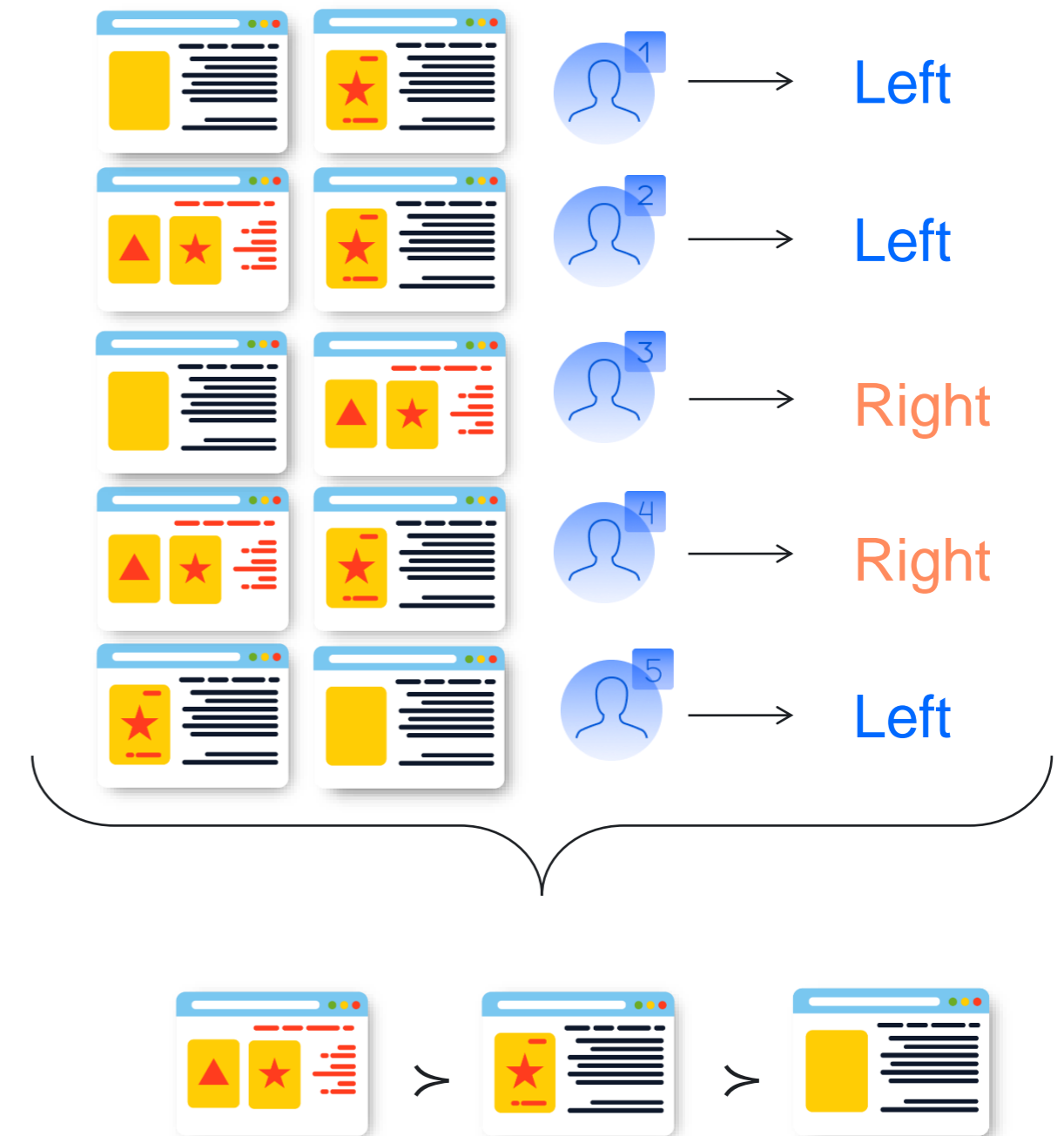
# Formalization

## 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**  $\pi$  over items  $D \rightarrow \{1, \dots, N\}$  based on answers in  $P$



# 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  $d_j \in D$

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

# Bradley and Terry Model: Example

| Performer | Task  | Left     | Right    |
|-----------|-------|----------|----------|
| $w_1$     | $t_1$ | <b>a</b> | b        |
| $w_1$     | $t_2$ | <b>b</b> | c        |
| $w_1$     | $t_3$ | c        | <b>a</b> |
| $w_2$     | $t_1$ | <b>a</b> | b        |
| $w_2$     | $t_2$ | <b>b</b> | c        |
| $w_2$     | $t_3$ | <b>c</b> | a        |

| Item | Score |
|------|-------|
| a    | 0.592 |
| b    | 0.278 |
| c    | 0.130 |

The model assumes that all performers are equally good and truthful!

# NoisyBT Model: Parameterization of Performers

$w_k$    $\longleftrightarrow$  “reliability”  $\gamma_k$  and “bias”  $q_k$

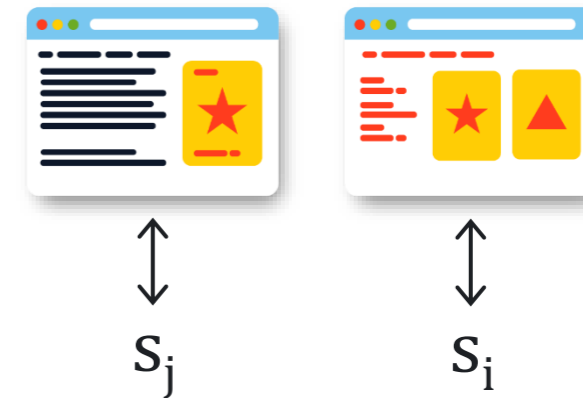
- ▶ The probability that  $w$  reads task is

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

- ▶ 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

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

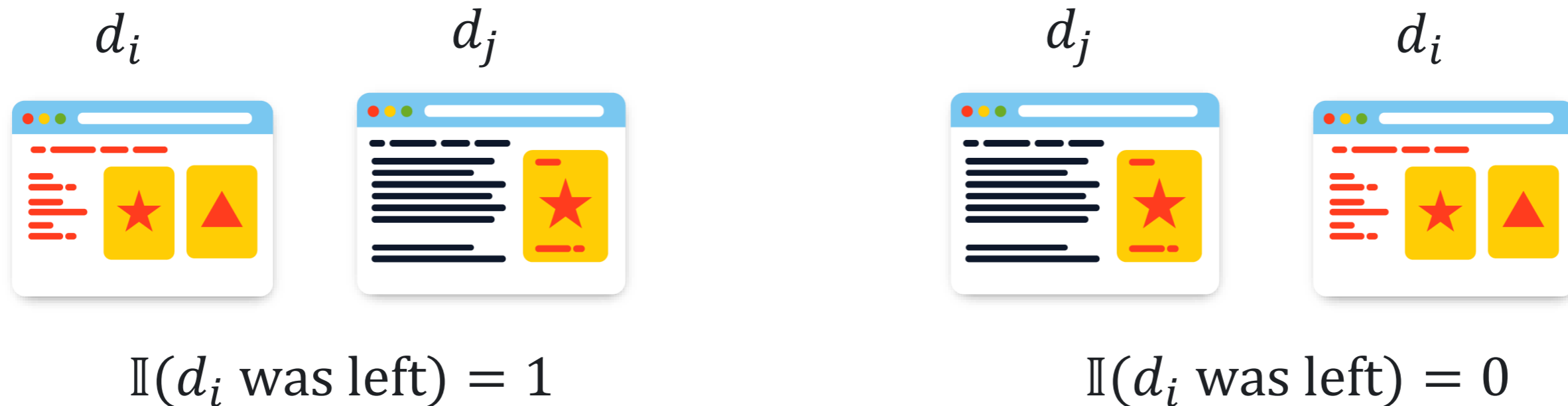
Probability to choose **Left** if answers randomly

# NoisyBT: Answer Likelihood

The likelihood of  $i \succ_k j$  is

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

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



# NoisyBT: Parameter Estimation

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**  $\pi$  over items, **sort** items according to their **scores**



# NoisyBT: Example

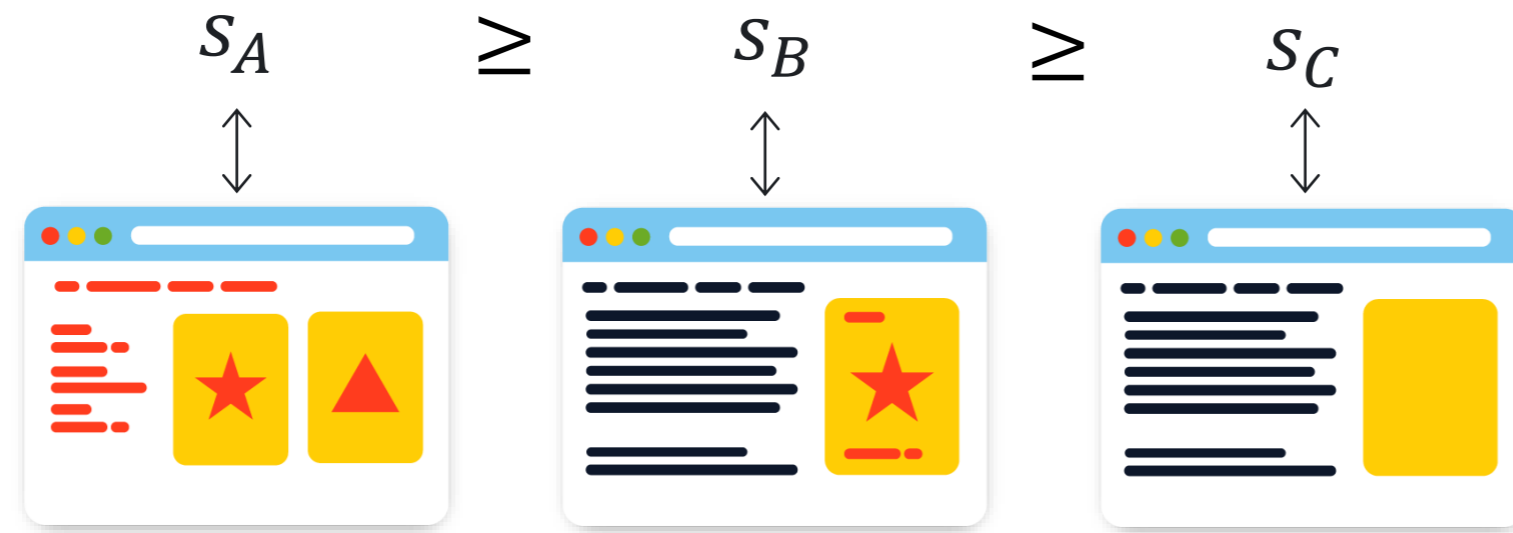
| Performer | Task  | Left     | Right    |
|-----------|-------|----------|----------|
| $w_1$     | $t_1$ | <b>a</b> | b        |
| $w_1$     | $t_2$ | <b>b</b> | c        |
| $w_1$     | $t_3$ | c        | <b>a</b> |
| $w_2$     | $t_1$ | <b>a</b> | b        |
| $w_2$     | $t_2$ | <b>b</b> | c        |
| $w_2$     | $t_3$ | <b>c</b> | a        |

| Item | Score |
|------|-------|
| a    | 1.000 |
| b    | 0.547 |
| c    | 0.000 |

| Performer | Bias  | Skill |
|-----------|-------|-------|
| $w_1$     | 0.633 | 0.656 |
| $w_2$     | 1.000 | 0.000 |

# Summary about NoisyBT

- ▶ Latent scores models for ranking from pairwise comparisons:



- ▶ To reduce bias from unreliable answers parameterize workers



**Demo**

# Demo

- ▶ We will learn how to aggregate your results using the Bradley-Terry model **right now**
- ▶ We will show you a **live demo**
- ▶ This demo will use the annotated **data from the practice session** that will be aggregated to provide the final rankings
- ▶ Please use a blank **Jupyter Notebook**, e.g., <https://colab.research.google.com/> or the local one

Crowd-Kit: <https://github.com/Toloka/crowd-kit>

**Crowd-Kit** is an open source **Python** library that implements methods for quality control in crowdsourcing.

- ▶ It is platform-agnostic, so can be used with any crowdsourcing platform
- ▶ It includes efficient implementations of classic and state-of-the-art methods for quality control
- ▶ It provides a simple API for using them in your application and is available on **PyPI**: `pip install crowd-kit`

# NoisyBT with Crowd-Kit: Input

```
import pandas as pd # pip install pandas
from crowdkit.aggregation import NoisyBradleyTerry # pip install -U crowd-kit

# In this example we will use the annotation results in the Toloka TSV format,
# but Crowd-Kit is platform-agnostic and it can handle any other format
df = pd.read_csv('assignments.tsv', sep='\t', dtype=str)
df.drop(df[~df['OUTPUT:result'].isin({'L', 'R'})].index, inplace=True)

# We need to reorganize our data frame to contain the following columns:
# query, performer, left, right, label
```

# NoisyBT with Crowd-Kit: Rename

```
df['performer'] = df['ASSIGNMENT:worker_id']
```

```
df['left'] = list(zip(df['INPUT:query'], df['INPUT:link_left']))
```

```
df['right'] = list(zip(df['INPUT:query'], df['INPUT:link_right']))
```

```
df.loc[df['OUTPUT:result'] == 'L', 'label'] = \  
df.loc[df['OUTPUT:result'] == 'L', 'left']
```

```
df.loc[df['OUTPUT:result'] == 'R', 'label'] = \  
df.loc[df['OUTPUT:result'] == 'R', 'right']
```

# NoisyBT with Crowd-Kit: Apply and Output

```
bt = NoisyBradleyTerry(n_iter=5000)
result = bt.fit_predict(df)

index = pd.MultiIndex.from_tuples(result.index)
df_result = pd.DataFrame(result, columns=['score'], index=index)
df_result['query'] = result.index.str[0]
df_result['url'] = result.index.str[1]
df_result.reset_index(drop=True, inplace=True)
df_result.sort_values(['query', 'score'], ascending=[True, False],
inplace=True)
df_result[['query', 'url', 'score']]
df_result.to_csv('aggregated.tsv', sep='\t', index=False,
                columns=['query', 'url', 'score'])
```



**Get Ready!**

# Conclusion

# Conclusion

- ▶ Human-in-the-Loop improves Web ranking by gathering **high-quality pairwise comparisons**
- ▶ Aggregated rankings can be used for evaluating your service against data from the **real users**
- ▶ **Crowd-Kit** allows aggregating human judgements in a simple way: <https://github.com/Toloka/crowd-kit>  
or `pip install crowd-kit`

# Thank you!

# Questions?

**Dmitry Ustalov**

Analyst/Software Developer



[dustalov@yandex-team.ru](mailto:dustalov@yandex-team.ru)



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