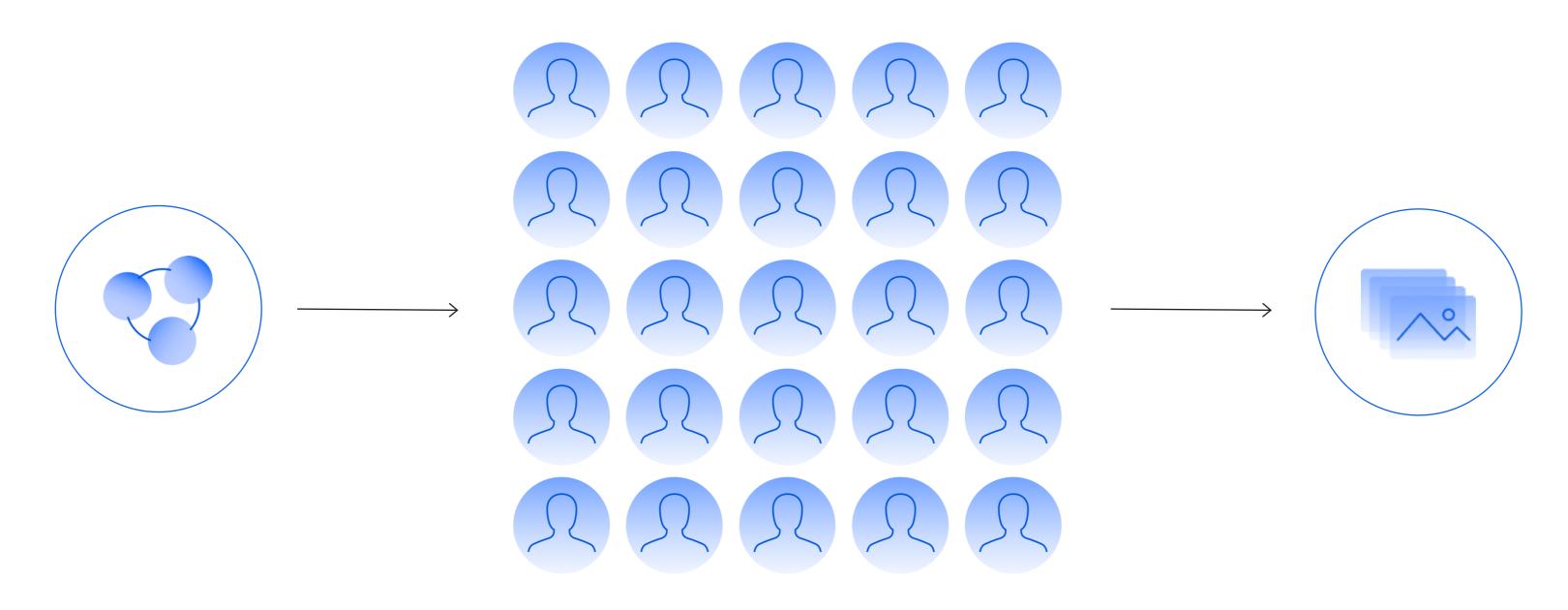Toloka

# Practice of Efficient Data Collection via Crowdsourcing: Aggregation, Incremental Relabelling, and Pricing

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

# Introduction

Olga Megorskaya, CEO, Toloka

# Crowdsourcing: specific way to design a business process



A big task                    Cloud of performers                    Result
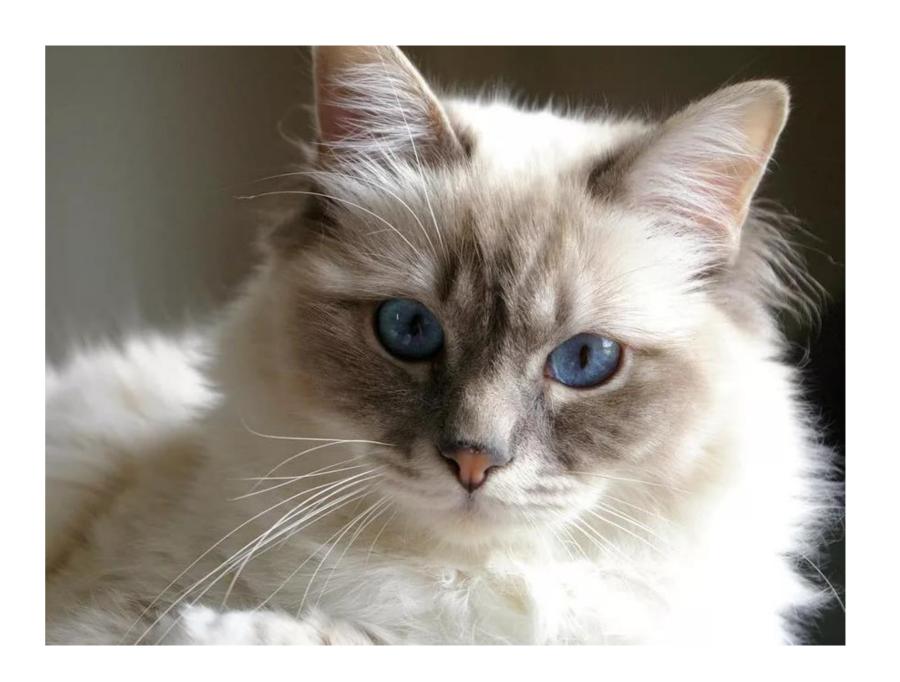
# Crowdsourcing applications: examples

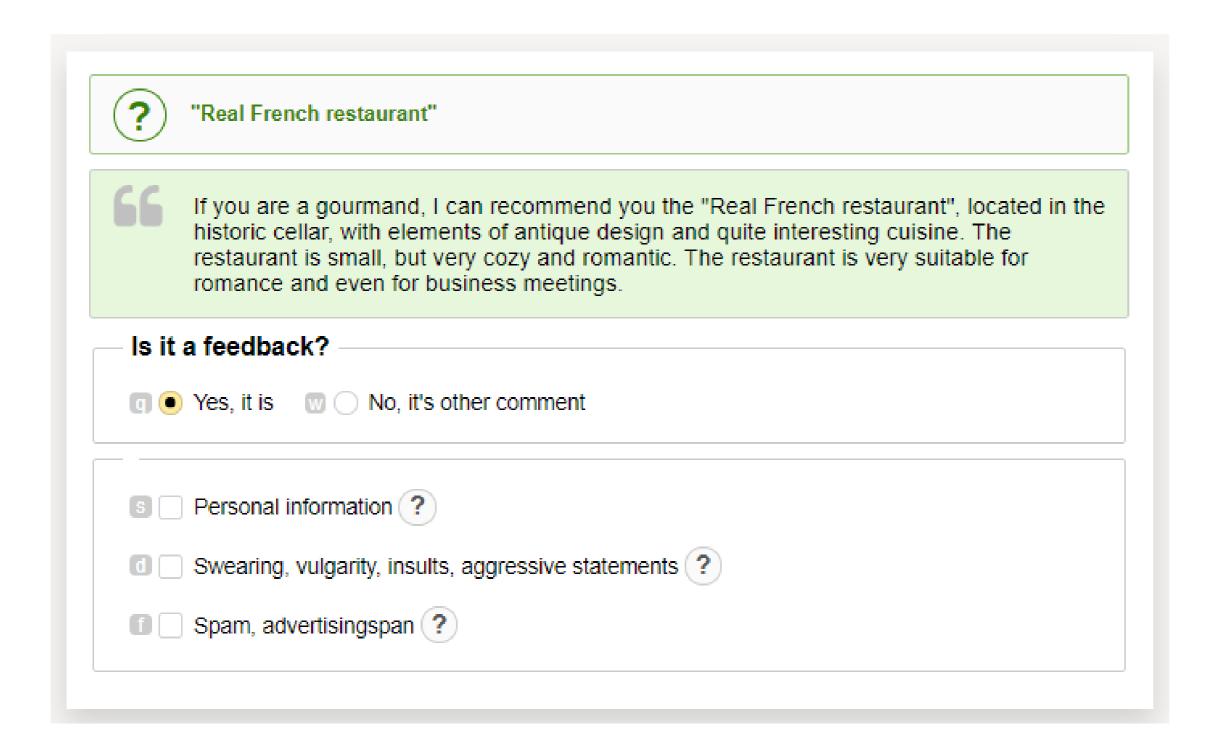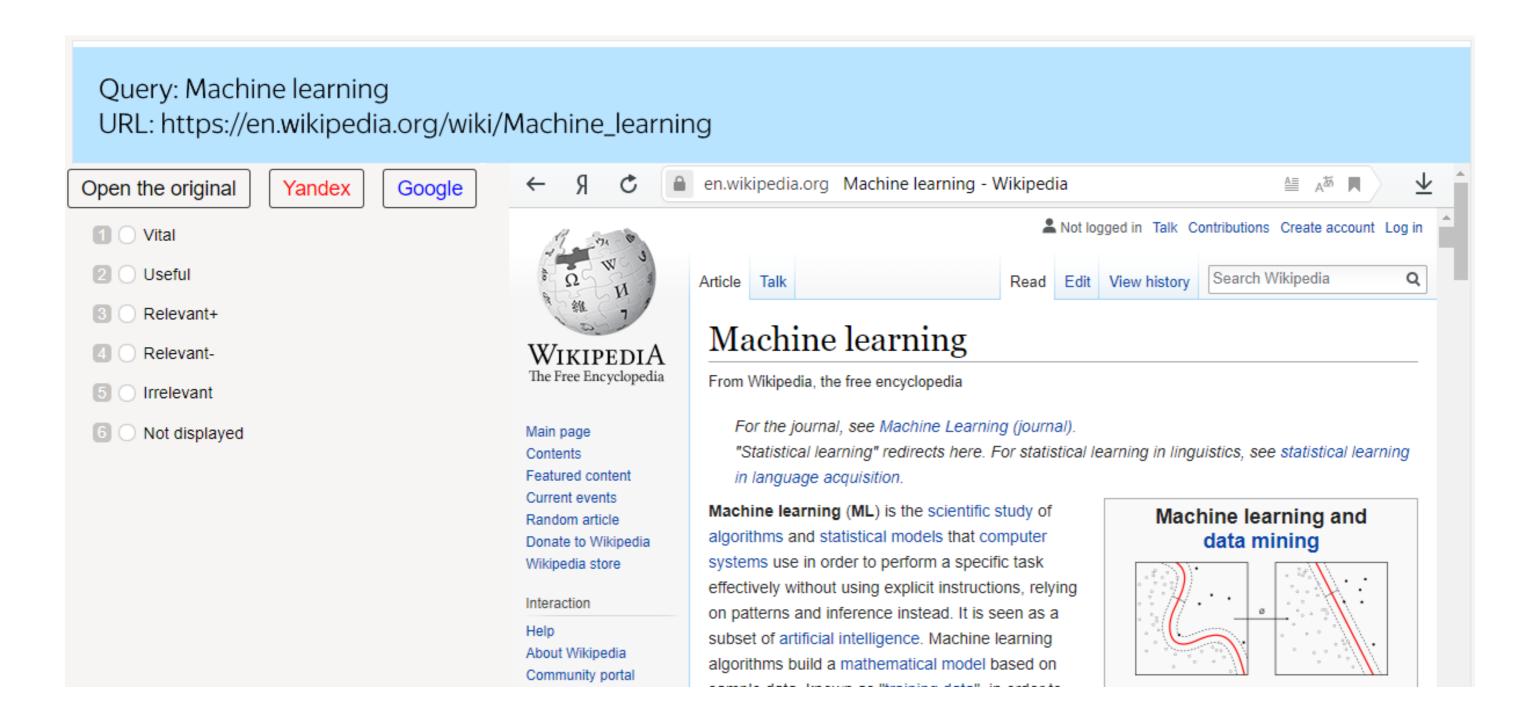| Task type | Where is used |
|---|---|
| Information assessment | Ranking of search results |
| Content categorization | Text and media moderation, data cleaning and filtering |
| Content annotation | Metadata tagging |
| Pairwise comparison | Offline evaluation, media duplication check |
| Object segmentation, including 3D | Image recognition for self-driving car |
| Audio and video transcription | Speech recognition for voice-controlled virtual assistant |
| Spatial crowdsourcing | Verify business information and office hours |

# Example: binary classification
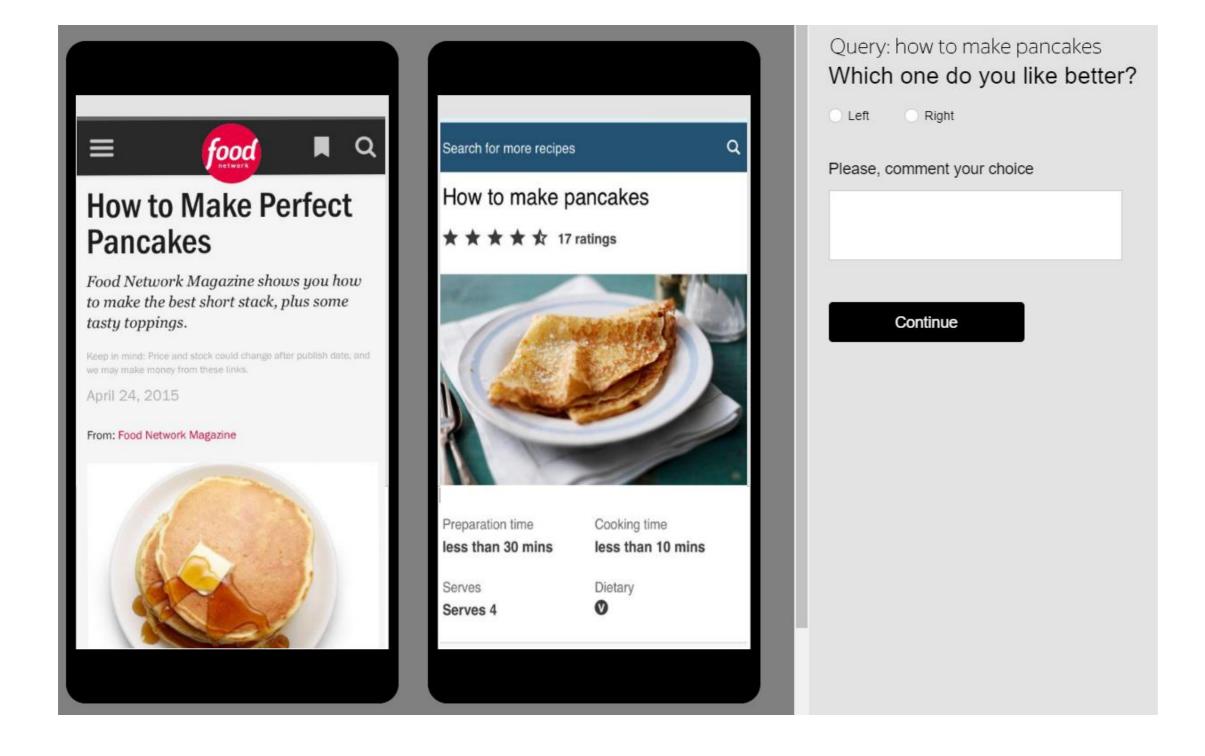
Is this cat white?

Yes

No

# Example: multi classification



**?** "Real French restaurant"

> If you are a gourmand, I can recommend you the "Real French restaurant", located in the historic cellar, with elements of antique design and quite interesting cuisine. The restaurant is small, but very cozy and romantic. The restaurant is very suitable for romance and even for business meetings.

**Is it a feedback?**

q ● Yes, it is    w ○ No, it's other comment

s ☐ Personal information **?**

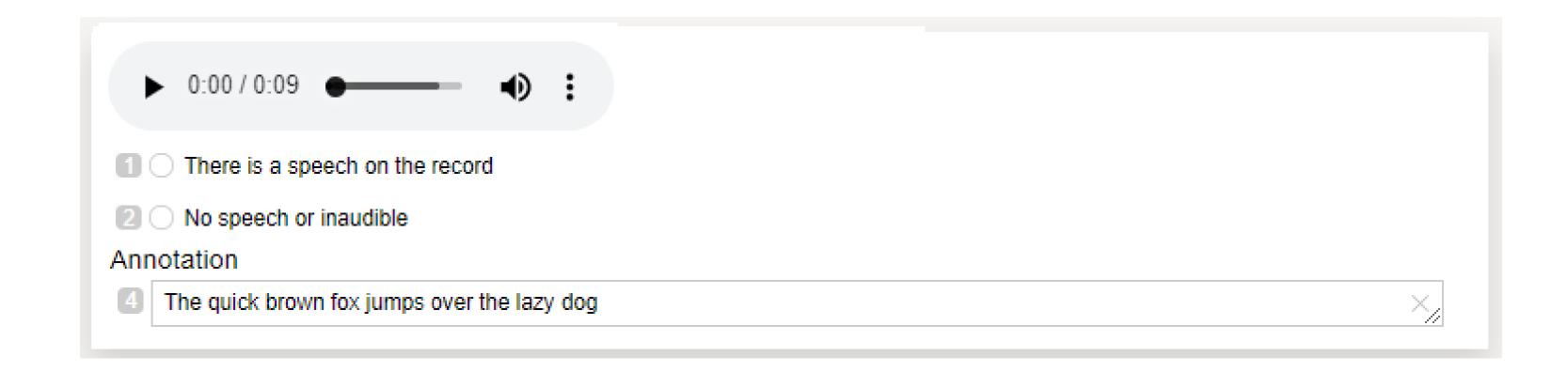d ☐ Swearing, vulgarity, insults, aggressive statements **?**

f ☐ Spam, advertisingspan **?**
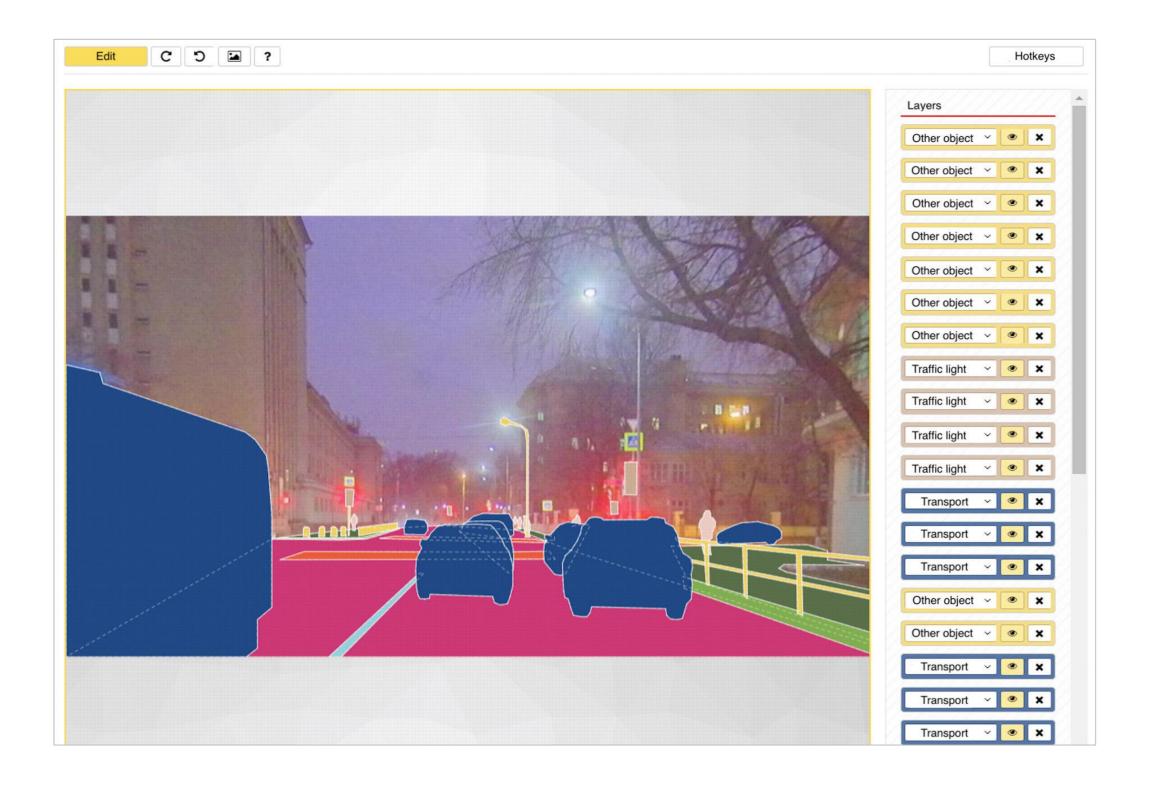
# Example: multi classification
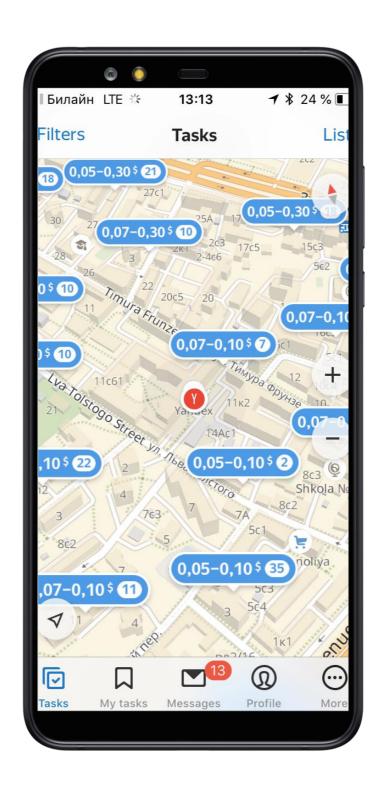# with ordered labels

# Examples: pairwise comparison

# Examples: transcription with textual answers

# Examples: object segmentation

# Examples: spatial crowdsourcing

# A crowdsourcing platform: two-sided market



Performers           Platform           Requesters

# Crowdsourcing platforms: examples

- ► Amazon Mechanical Turk

- ► Toloka

- ► Microworkers

- ► Gigwalk

- ► ClickWorker

- ► CloudFactory

- ► Figure Eight

- ► CrowdSource

- ► DefinedCrowd

- ► …

# Pros of crowdsourcing platforms

24/7

Variety
of skilled
performers

Vast
region
coverage

Ongoing
processes

# Crowdsourcing growth: our experience

Active performers in Toloka



| 2014 | 2015 | 2016 | 2017 | 2018 | 2019 |
|------|------|------|------|------|------|
| 9K | 120K | 270K | 570K | 1.1M | 2.2M |

# Crowdsourcing growth: our experience

Different projects in Toloka

| Year | Value |
|------|-------|
| 2014 | 57 |
| 2015 | 443 |
| 2016 | 1214 |
| 2017 | 1612 |
| 2018 | 2123 |
| 2019 | 4055 |

# Everyday on Toloka

500+
different
projects

36K+
performers

12M+
tasks

# Toloka: real-life cases

| Case | Tasks | Done in | Cost |
|------|-------|---------|------|
| Side-by-side object comparison | 1,000 tasks | 10 min | $2.4 |
| Object classification | 1,000 photos | 15 min | $1.2 |
| Object segmentation | About 1,000 objects in 100 photos | 6 h | $3.6 |
| Phrase generation for a chatbot | 500 phrases for the same topic | 15 min | $1 |
| Audio transcription | 100 recordings 25 minute long | 20 min | $6 |
| Video ranking | 10,000 videos | 2 h | $10 |

# Tutorial overview

# Why this tutorial?
**Practice**

# Part I: 30 min

## Main components of data collection via crowdsourcing

▶ Decomposition for effective pipeline

▶ Task instruction & interface: best practices

▶ Quality control techniques



Olga Megorskaya

CEO, Toloka

# Part II: 25 min

## Analysis of label collection projects to be done (practical session)

► Dataset and required labels

► Discussion: how to collect labels?

► Data labelling pipeline for implementation

Daria Baidakova

Project Manager, Toloka

# Part III: 10 min

## Introduction to the crowdsourcing platform Yandex.Toloka for requesters

- ► Main types of instances
- ► Project: creation & configuration
- ► Pool: creation & configuration
- ► Tasks: uploading & golden set creation
- ► Statistics in flight and download of results

Evfrosiniya Zerminova

Head of Data Analysis and Research Group, Toloka

# Part IV: 60 min

## Setting up and running label collection projects (practical session)

You

► Create

► Configure

► Run on real performers

Data labelling projects in real-time

Daria Baidakova

Project Manager,
Toloka

# Part V: 35 min

## Interface & quality control

▶ Detailed examination of quality control techniques

▶ Comprehensive overview of best practices for creating a functional interface

Alexey Drutsa

Head of Efficiency and Growth Division, Toloka

# Part VI: 25 min

## Theory on Aggregation

► Multiclass labels

► Pairwise comparisons



Valentina Fedorova

Researcher, Toloka

# Part VII: 90 min

## Setting up and running label collection projects cont. (practical session)

You

► Create

► Configure

► Run on real performers

Data labelling projects in real-time

Daria Baidakova

Project Manager,
Toloka

# Part VIII: 20 min

## Theory on efficient incremental relabelling and pricing

▶ Incremental relabelling

▶ Performance-based pricing

Valentina Fedorova

Researcher, Toloka

# Part IX: 10 min

## Discussion of results from the projects & conclusions

► Results of your projects

► Extensions to work on after tutorial



Alexey Drutsa

Head of Efficiency and Growth Division, Toloka

# Tutorial outline

**Introduction:**
**20 min**

**Part II: 25 min**
Brainstorming pipeline

**Lunch break:**
**90 min**

**Coffee break:**
**30 min**

**Part I: 40 min**
Main Components

**Part III: 10 min**
Introduction to Crowd Platform

**Part V: 35 min**
Interface & Quality control

**Part VII: 60 min**
Set & Run Projects cont.

**Coffee break:**
**30 min**

**Part IV: 85 min**
Set & Run Projects

**Part VI: 25 min**
Theory on Aggregation

**Part VIII: 20 min**
Incremental relabeling and pricing

**Part IX: 10 min**
Results & Conclusions