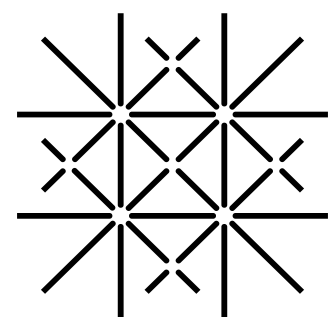# Intrinsic Self-Supervision for Data Quality Audits

**Fabian Gröger, Simone Lionetti, Philippe Gottfrois, Alvaro Gonzalez-Jimenez, Ludovic Amruthalingam, Matthew Groh, Alexander A. Navarini, Marc Pouly**
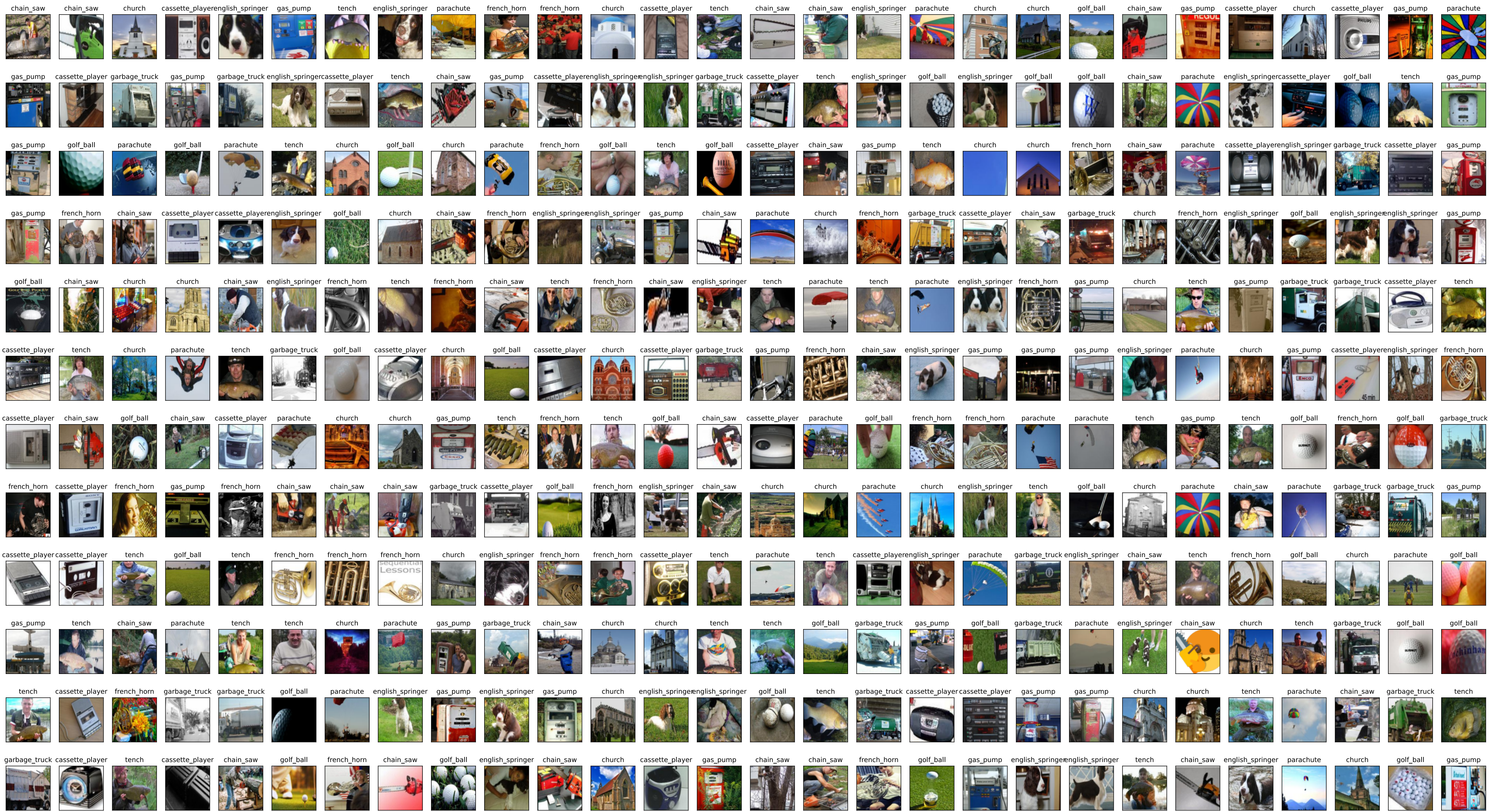
University of Basel

**HSLU** Lucerne University of Applied Sciences and Arts

Northwestern University

Off-topic samples    Near duplicates    Label errors

# Everyone who is doing ML knows …

… training AND evaluation data can be messy.

… noise during evaluating leads to inconsistent performance estimates.

# Everyone who is doing ML knows …

… training AND evaluation data can be messy.

… noise during evaluating leads to inconsistent performance estimates.

# BUT, data cleaning …

… can be very time-consuming and labor intensive.
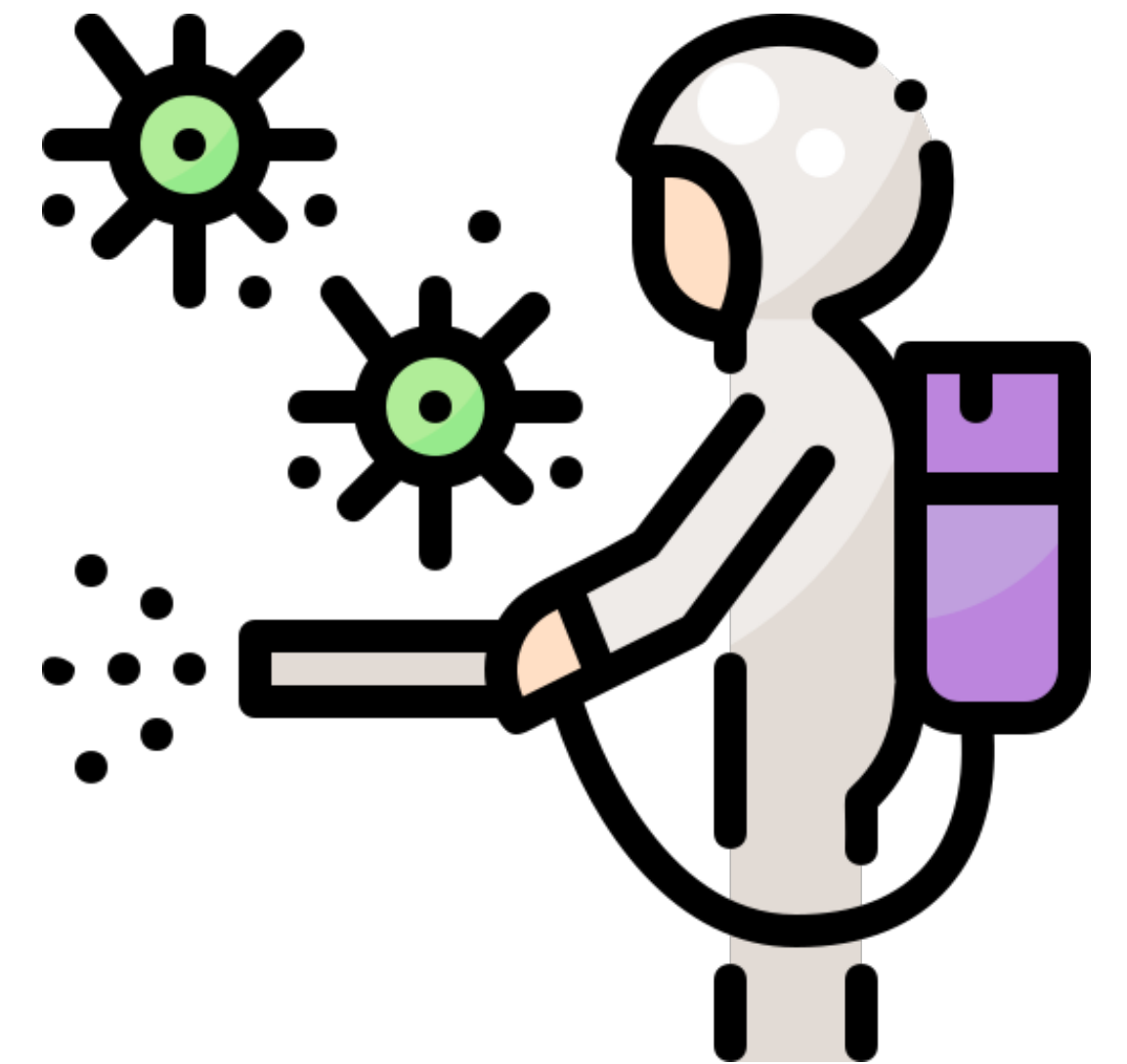
… is the least enjoyable task for many practitioners*.
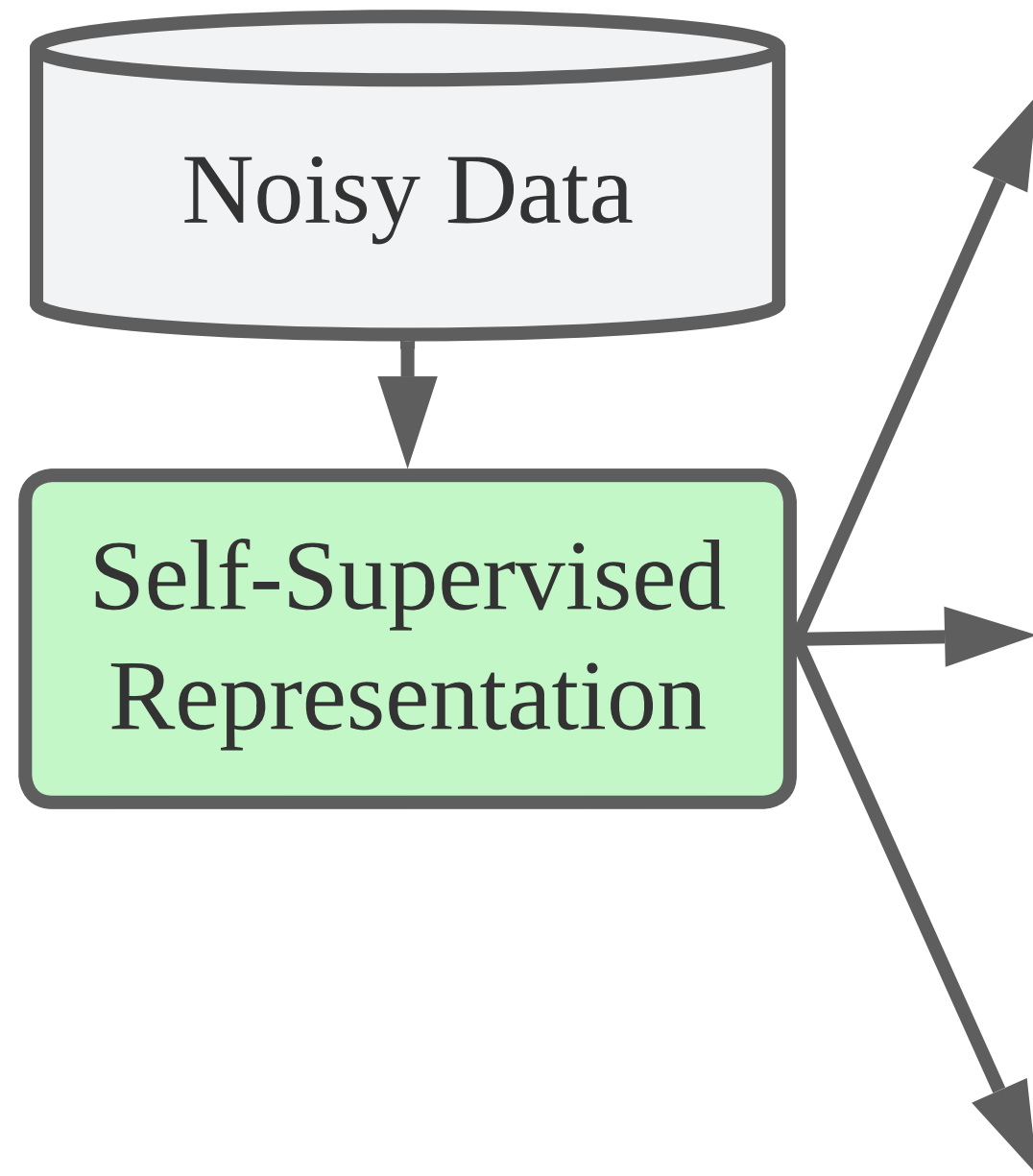
# Goals of this project 🏆

1. Reliably **detect data quality issues**, such as off-topic images, near duplicates, and label errors, in image datasets without introducing significant biases.

2. **Reduce the time needed** for detecting and confirming data quality issues.

3. Investigate the **influence of data quality issues** on model training and evaluation.

# Our findings

- Self-supervised learned (SSL) **representations can be exploited to find data quality issues**.

- **Context-aware SSL** representations can capture the dataset context with minimal bias.

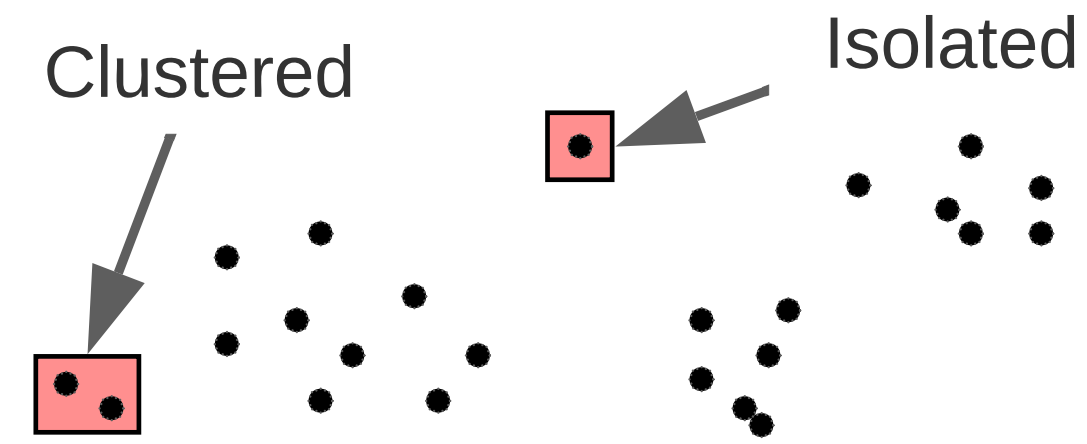- Combination of SSL **representations and distance-based indicators** effectively finds quality issues.
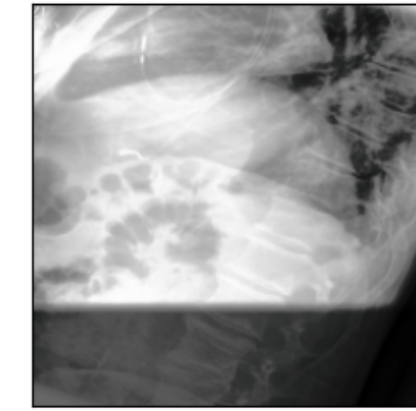
**SelfClean**

Noisy Data

Self-Supervised Representation

**Off-topic Samples**
Agglomerative clustering

Clustered

Isolated

**Near Duplicates**
Pairwise distance

Approximate Duplicate

Exact Duplicate

**Label Errors**
Intra-/extra- class distance ratio

Label Errors

ImageNet    CheXpert    Fitzpatrick17k

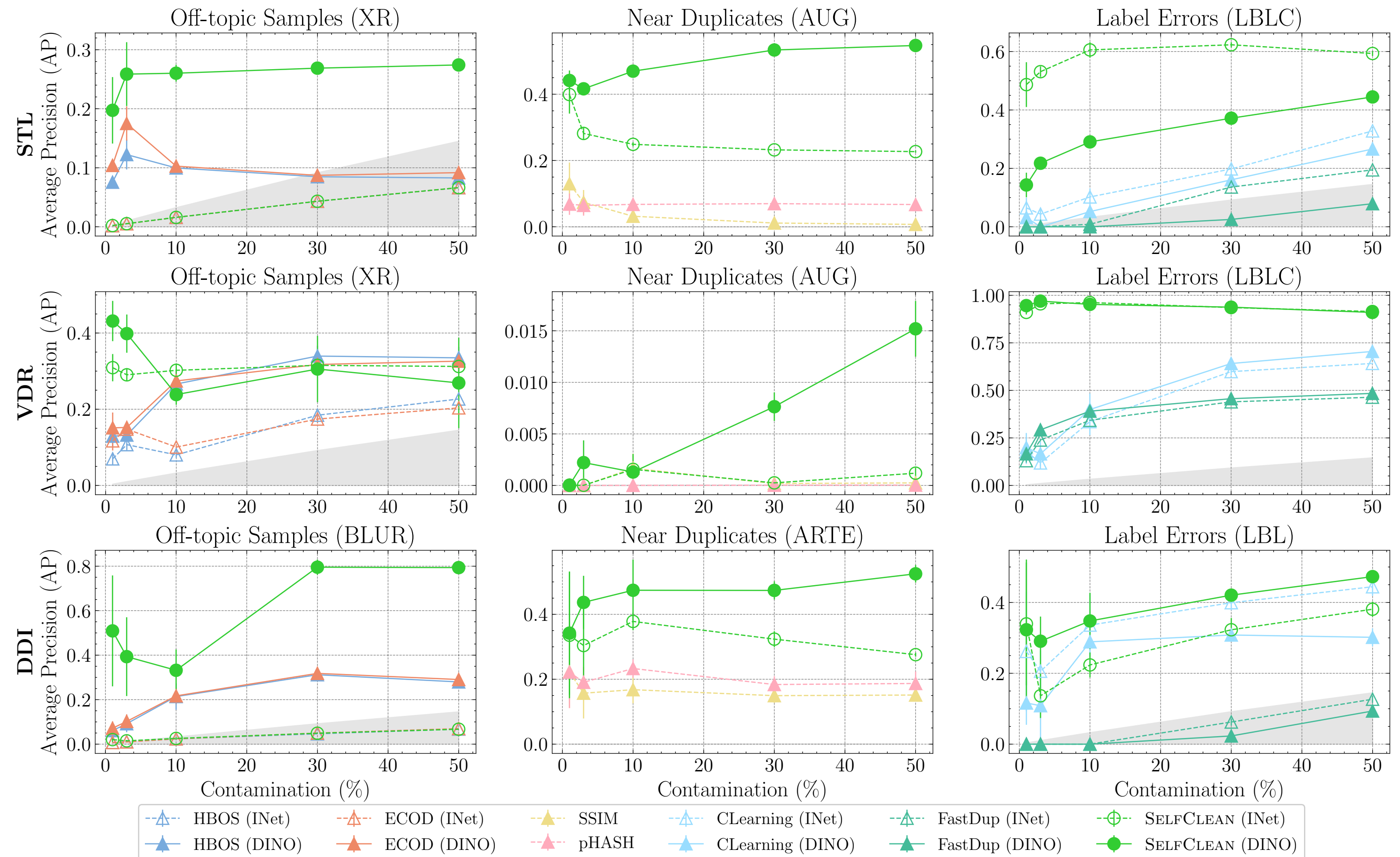Kite    Atelectasis: positive    Benign epidermal

# Results I

- Evaluation on both synthetic and natural contamination showed a **significant improvement** compared to current solutions.
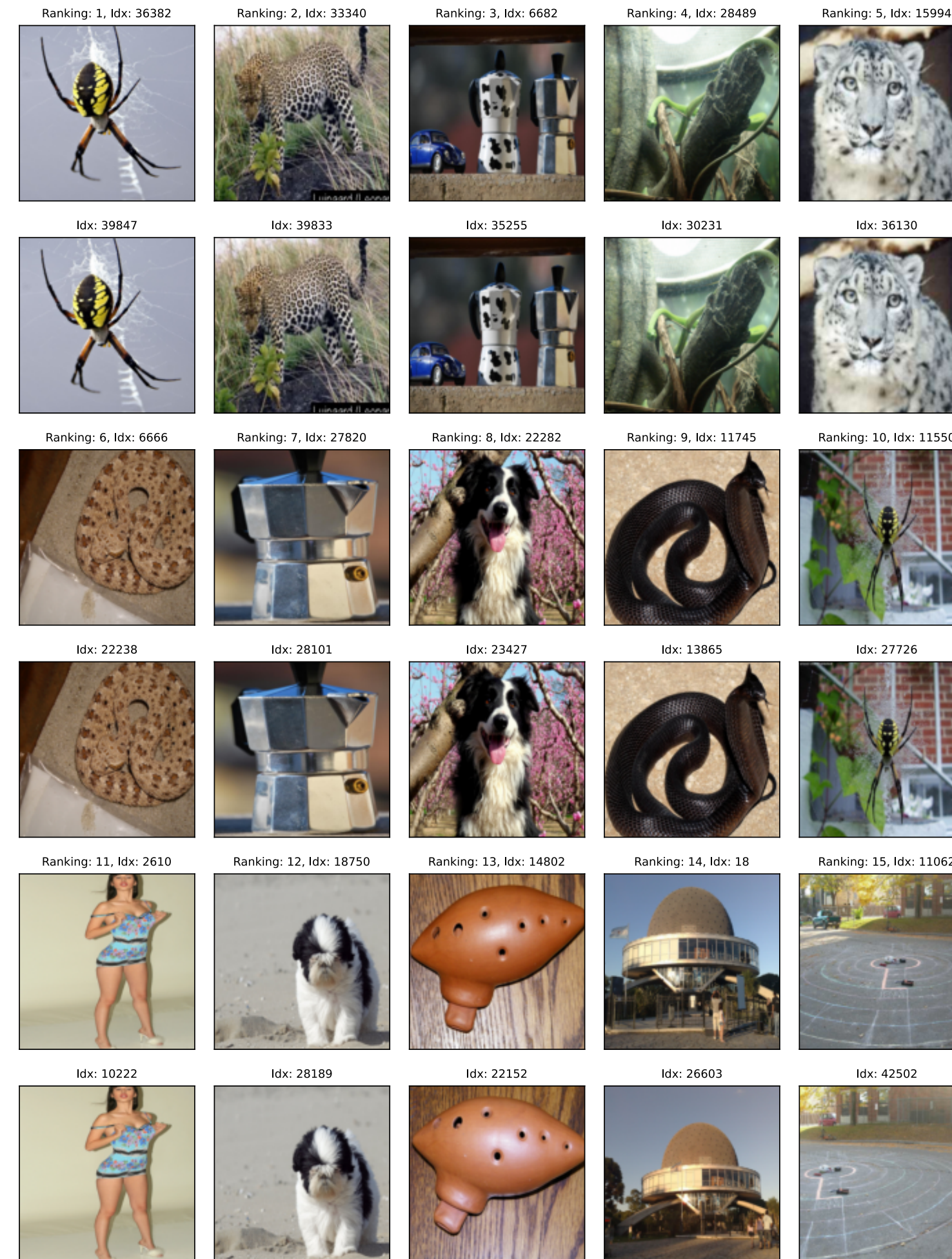


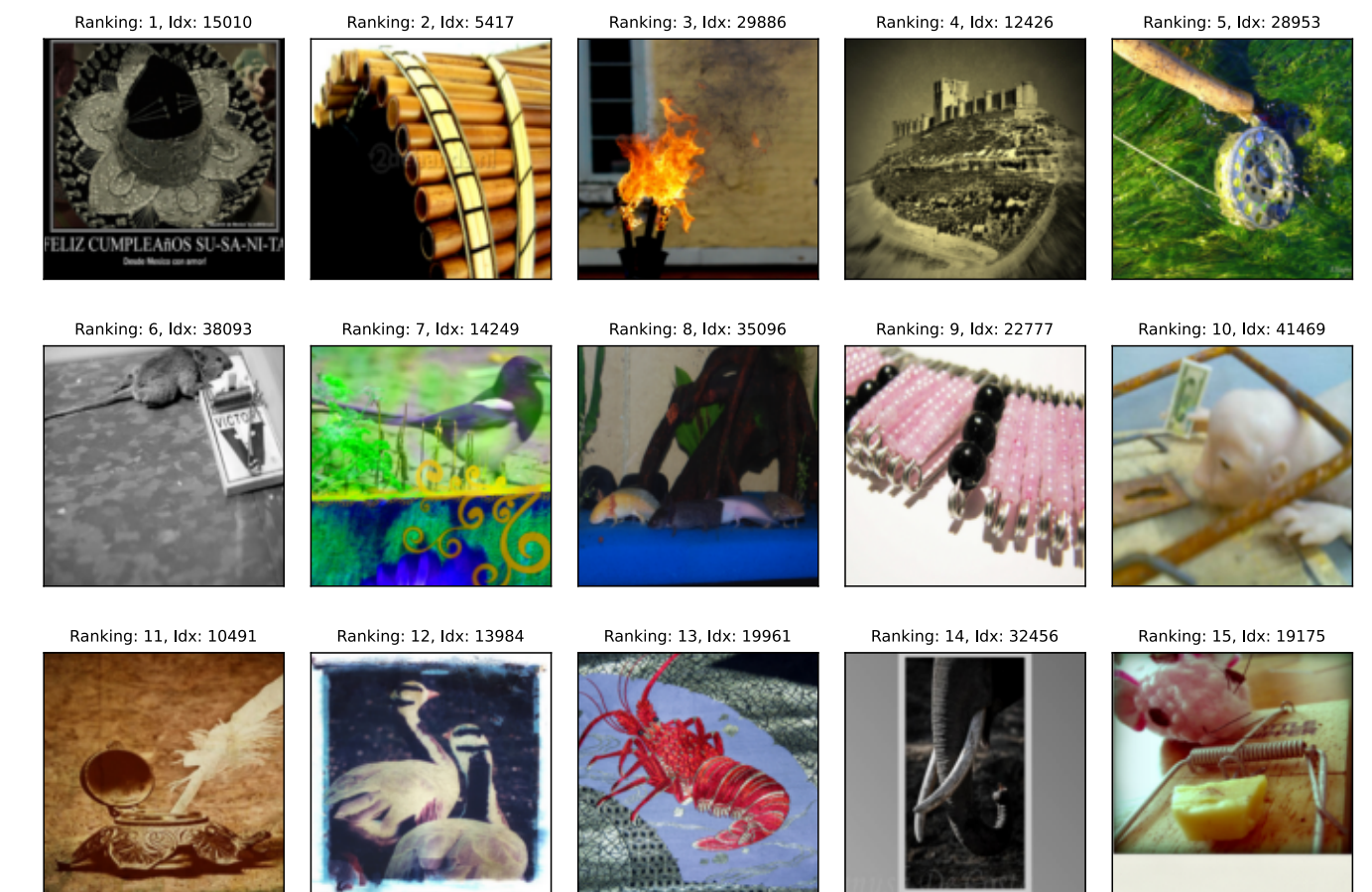Synthetic evaluation results.

# Results II

- Applied to multiple image benchmarks, we identify up to **16% of issues**, and confirm an **improvement in evaluation reliability** upon cleaning.
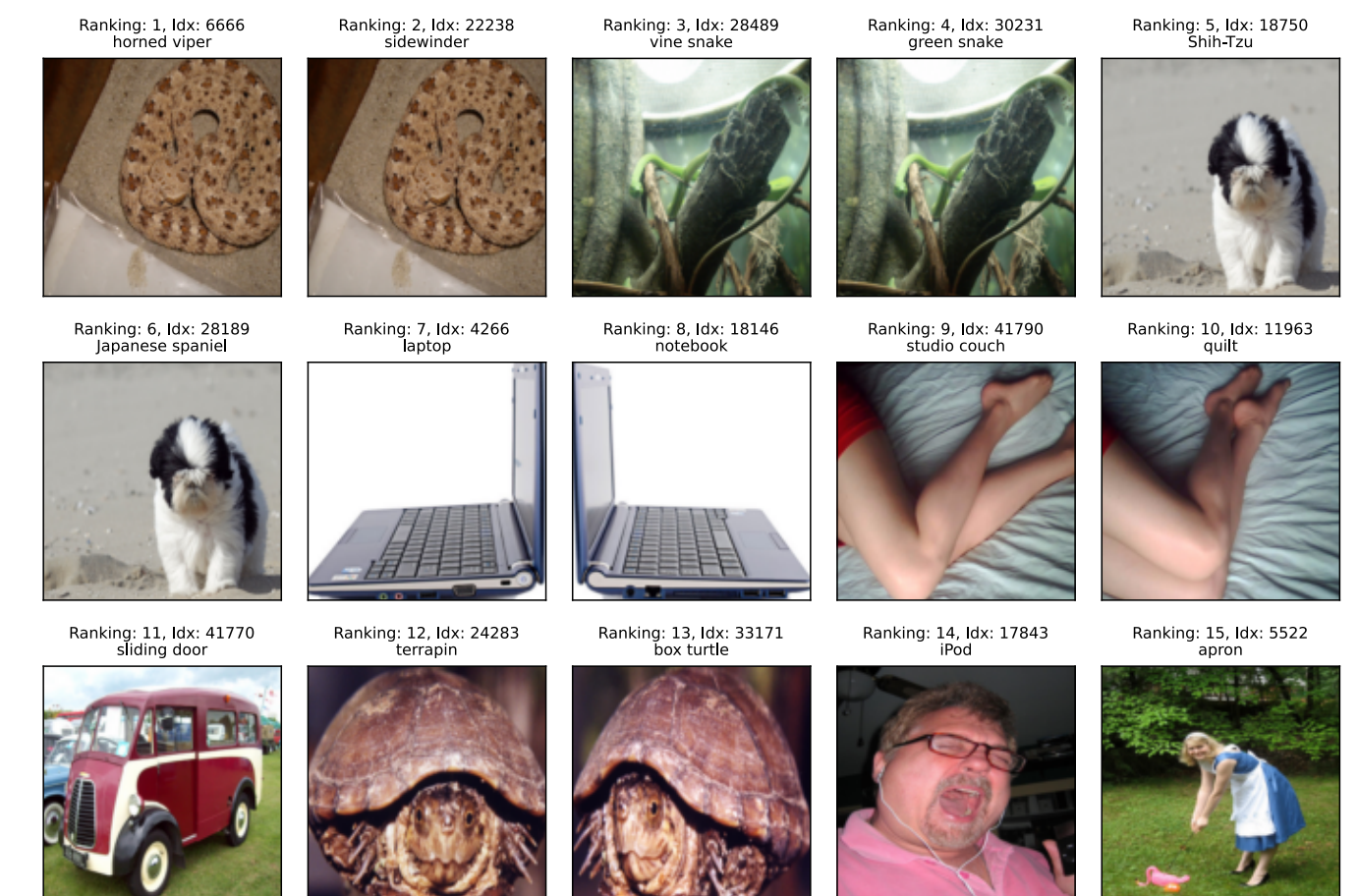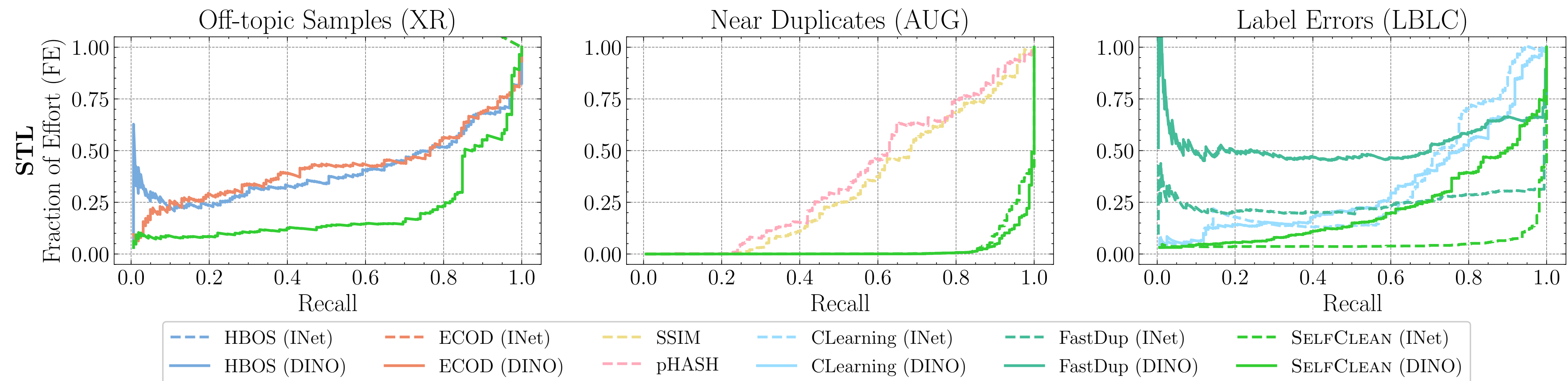


(a) Near duplicates

(b) Off-topic samples

(c) Label errors

Analysis of ImageNet-1k.

# Results III

- For a typical dataset SelfClean can **reduce the inspection effort by a factor between 5 and 50**.



Analysis of the inspection effort saved.

# Intrinsic Self-Supervision for Data Quality Audits

**Project website: selfclean.github.io/**